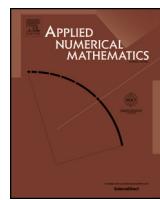




Contents lists available at ScienceDirect



Diagonally implicit Runge–Kutta methods for stiff ODEs

Christopher A. Kennedy ^{a,*}, Mark H. Carpenter ^b^a Palo Alto, CA, USA^b Computational Modeling and Simulation Branch, NASA Langley Research Center, Hampton, VA 23681, USA

ARTICLE INFO

Article history:

Received 15 June 2019

Received in revised form 7 July 2019

Accepted 8 July 2019

Available online 11 July 2019

Keywords:

Embedded method

Simplifying assumptions

Singular perturbation equations

L-stable

L-stability

ABSTRACT

Based principally on a recent review of diagonally implicit Runge–Kutta (DIRK) methods applied to stiff first-order ordinary differential equations (ODEs) by the present authors, several nearly optimal, general purpose, DIRK-type methods are presented. Methods given range from third- to sixth-order in four- to nine-stages. All are both stiffly-accurate and L-stable and many are internally L-stable on stages where this is possible. Focus is placed on stage-order two methods. To facilitate step-size control via error estimation, an embedded method is included with each method listed. As the present methods must make no concessions to an explicit method as part of an additive Runge–Kutta method, the present ESDIRKs are generally better than their IMEX analogs. Five new explicit first-stage, singly diagonally-implicit Runge–Kutta (ESDIRK) methods are presented based on lessons learned from the review. To the best of our knowledge, the sixth-order, L-stable, stage-order two, 6(5)-pair, ESDIRK included herein is the only one of its kind.

© 2019 IMACS. Published by Elsevier B.V. All rights reserved.

1. Introduction

The diagonally implicit Runge–Kutta (DIRK) family of methods is possibly the most widely used implicit Runge–Kutta (IRK) method in practical applications involving stiff, first-order, ordinary differential equations (ODEs) for initial value problems (IVPs) due to their relative ease of implementation. They are characterized by a lower triangular A -matrix with at least one nonzero diagonal entry and are sometimes referred to as semi-implicit or semi-explicit Runge–Kutta methods. This structure permits solving for each stage individually rather than all stages simultaneously.

Since the first effectively multistage DIRK-type method was introduced by Butcher [5], much has been learned about IRK methods in general, and DIRK-type methods in particular. Descriptions of this class of implicit Runge–Kutta methods may be found in the books by Burrage [4], Butcher [6], Dekker and Verwer [7] and Hairer & Wanner [9] as well as an older review paper by Alexander [1]. Beyond use as stand-alone methods, it is common to see DIRK-type methods employed as one part of a partitioned Runge–Kutta method. An implicit-explicit (IMEX) or additive partitioning of Runge–Kutta methods, pioneered by Griepentrog [8] and others, tailors the integration method to the characteristics of the equations being integrated. Examples of these IMEX methods may be found in the papers by Ascher, Ruuth and Spiteri [2], Kennedy and Carpenter [13,15] and more recently, Izzo and Jackiewicz [12].

DIRK methods have been the subject of a recent review [14] by the present authors. This paper is based on that review. In the review, many aspects of DIRK-type methods are discussed which are relevant to proper design and application of such methods. These include the general structure, order conditions, simplifying assumptions, error, linear stability, nonlinear sta-

* Corresponding author.

E-mail address: cakenne@gmail.com (C.A. Kennedy).

bility, internal stability, dense output, conservation, symplecticity, symmetry, dissipation and dispersion accuracy, memory economization [11], regularity, boundary and smoothness order reduction, efficiency [3], solvability, implementation, step-size control, iteration control, stage-value predictors, discontinuities and existing software. A comprehensive bibliography on DIRK methods is also included in the review with 470 citations. Hence, references in the present paper will be limited and focused on more recent contributions to the literature.

Several conclusions were reached in the review. Among them were that if general DIRK methods are intended to solve stiff equations, methods should possess a stage-order of two, L-stability, strong damping of stiff eigenvalues on the internal stages and minimized leading-order error. Such methods, along with error estimators, dense output and stage-value predictors form the basis of a robust and efficient integrator.

The present paper focuses on methods which are at least third-order accurate and introduces five new ESDIRK methods based on lessons learned from [14]. As the present methods must make no concessions to an explicit method as part of an additive Runge–Kutta method, the present ESDIRKs are generally better than their IMEX analogs [15]. Dense output and stage-value predictors are not included with these methods. Appendix A lists the method coefficients for five new ESDIRK methods.

2. Background

DIRK-type methods are used to solve ODEs of the form

$$\frac{dU}{dt} = F(t, U(t)), \quad U(a) = U_0, \quad t \in [a, b] \quad (1)$$

and are applied over s -stages as

$$\left. \begin{aligned} F_i &= F(t_i, U_i), & U_i &= U^{[n]} + (\Delta t) \sum_{j=1}^s a_{ij} F_j, & t_i &= t^{[n]} + c_i \Delta t, \\ U^{[n+1]} &= U^{[n]} + (\Delta t) \sum_{i=1}^s b_i F_i, & \hat{U}^{[n+1]} &= U^{[n]} + (\Delta t) \sum_{i=1}^s \hat{b}_i F_i, \end{aligned} \right\} \quad (2)$$

where $i = 1, 2, \dots, s$, $F_i = F_i^{[n]} = F(U_i, t^{[n]} + c_i \Delta t)$. Also, $\Delta t > 0$ is the step-size, $U^{[n]} \simeq U(t^{[n]})$ is the value of the U -vector at time step n , $U_i = U_i^{[n]} \simeq U(t^{[n]} + c_i \Delta t)$ is the value of the U -vector on the i th-stage, and $U^{[n+1]} \simeq U(t^{[n]} + \Delta t)$. Both $U^{[n]}$ and $U^{[n+1]}$ are of classical order p . The U -vector associated with the embedded scheme, $\hat{U}^{[n+1]}$, is of order $\hat{p} = p - 1$. This constitutes a (p, \hat{p}) pair. Each of the respective Runge–Kutta coefficients a_{ij} (stage weights), b_i (scheme weights), \hat{b}_i (embedded scheme weights), and c_i (abscissae or nodes), $i, j = 1, 2, \dots, s$ are real and are constrained, at a minimum, by certain order of accuracy and stability considerations.

For the stiffly-accurate ($a_{sj} = b_j$, $j = 1, 2, \dots, s$), stage-order two methods considered in this paper, ESDIRK methods are chosen. They are given by the general structure

$$\begin{array}{c|ccccccccc} 0 & 0 & 0 & 0 & \cdots & 0 & 0 & 0 & 0 \\ 2\gamma & \gamma & \gamma & 0 & \cdots & 0 & 0 & 0 & 0 \\ c_3 & a_{31} & a_{32} & \gamma & \ddots & 0 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \ddots & \ddots & \vdots & \vdots & \vdots \\ c_{s-2} & a_{s-2,1} & a_{s-2,2} & a_{s-2,3} & \ddots & \gamma & 0 & 0 & 0 \\ c_{s-1} & a_{s-1,1} & a_{s-1,2} & a_{s-1,3} & \cdots & a_{s-2,s-2} & \gamma & 0 & 0 \\ 1 & b_1 & b_2 & b_3 & \cdots & b_{s-1} & b_s & \gamma & \\ \hline & b_1 & b_2 & b_3 & \cdots & b_{s-2} & b_{s-1} & \gamma & \\ & \hat{b}_1 & \hat{b}_2 & \hat{b}_3 & \cdots & \hat{b}_{s-2} & \hat{b}_{s-1} & \hat{b}_s. & \end{array}$$

Some authors prefer to decompose $\mathbf{A} = a_{ij}$, $\mathbf{b} = b_i$ and $\mathbf{c} = c_i$, into

$$\mathbf{c} \left| \begin{array}{c|cc} \mathbf{A} & 0 & \mathbf{0}^T \\ \hline \mathbf{b}^T & \hat{\mathbf{c}} & \hat{\mathbf{A}} \end{array} \right. = \left[\begin{array}{c|cc} 0 & \mathbf{0}^T \\ \hline \mathbf{a} & \hat{\mathbf{A}} \end{array} \right], \quad \text{with} \quad \mathbf{A} = \left[\begin{array}{c|cc} 0 & \mathbf{0}^T \\ \hline \mathbf{a} & \hat{\mathbf{A}} \end{array} \right], \quad (3)$$

where $\mathbf{0}^T$ (composed of zeros), \mathbf{a}^T , $\hat{\mathbf{c}}^T$ and $\hat{\mathbf{b}}^T$ are vectors of length $(s-1)$, $\{0, \mathbf{a}^T\} = a_{11}^T$, $\{0, \hat{\mathbf{c}}^T\} = \mathbf{c}^T$, and $\{b_1, \hat{\mathbf{b}}^T\} = \mathbf{b}^T$ and $\hat{\mathbf{A}}$ is a square matrix of dimension $(s-1) \times (s-1)$. Hence, although \mathbf{A} is not invertible, $\hat{\mathbf{A}}$ is often invertible. The motivation for having an explicit first stage is primarily to allow stage-order two methods.

To identify certain schemes derived in this paper, they will be named ESDIRK $p(\hat{p})sS[q]X_x$, where p is the order of the main method, \hat{p} is the order of the embedded method, s is the number of stages, S is some stability characterization of the method, q is the stage-order of the method, X is used for any other important characteristic of the method, and x distinguishes between family members of some type of method.

3. Order conditions

Expressions for the equations of condition or truncation error coefficients, $\tau_j^{(p)}$, associated with the p th-order trees are of the form [6,10,9]

$$\tau_j^{(p)} = \frac{1}{\sigma} \Phi_j^{(p)} - \frac{\alpha}{p!} = \frac{1}{\sigma} \left(\Phi_j^{(p)} - \frac{1}{\gamma} \right), \quad p = 1, 2, \dots, j = 1, 2, \dots, a_p, \quad (4)$$

$$\Phi_j^{(p)} = \sum_{i=1}^s b_i \Phi_{i,j}^{(p)}, \quad \alpha \sigma \gamma = p! \quad (5)$$

where $a_p = \{1, 1, 2, 4, 9, 20, 48, 115, 286, 719\}$ for $p = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$. Also, α , ρ , γ are, respectively, the cardinality, the order and the symmetry of the particular tree which represents the order condition. The elementary weights, $\Phi_{i,j}^{(p)}$ and $\Phi_j^{(p)}$, are Runge–Kutta coefficient products and their sums, and j represents the index of the order condition, i.e. each order may have more than one order condition.

The local truncation error of an ERK is given by [16],

$$U(t^{[n]} + (\Delta t)) - U^{[n+1]} = \sum_{i=1}^{\infty} (\Delta t)^i \sum_{j=1}^{a_i} \tau_j^{(i)} \mathbf{F}_j^{(i)}(U^{[n]}), \quad (6)$$

where $\mathbf{F}_j^{(i)}(U^{[n]}) = \mathbf{F}_j^{(i)}$ are elementary differentials. A Runge–Kutta method is said to be of order p if the local truncation error satisfies

$$U^{[n]} - U(t^{[n]}) = \mathcal{O}(\Delta t)^{p+1}. \quad (7)$$

No order conditions are included in this study specific to differential-algebraic equations.

3.1. Simplifying assumptions

Simplifying assumptions are often made to facilitate the solution of Runge–Kutta order conditions, $\tau_j^{(p)}$, and possibly to enforce other desirable method characteristics such as higher stage-order. The three common ones are [6,9]

$$B(p) : \sum_{j=1}^s b_i c_i^{k-1} = \frac{1}{k}, \quad k = 1, \dots, p, \quad (8)$$

$$C(\eta) : \sum_{j=1}^s a_{ij} c_j^{k-1} = \frac{c_i^k}{k}, \quad k = 1, \dots, \eta, \quad (9)$$

$$D(\zeta) : \sum_{i=1}^s b_i c_i^{k-1} a_{ij} = \frac{b_j}{k} (1 - c_j^k), \quad k = 1, \dots, \zeta. \quad (10)$$

A fourth simplifying assumption, $E(\eta, \zeta)$ [7], is not needed for DIRK-type methods. The first is related to the quadrature conditions. Enforcing $B(p)$ forces $\tau_1^{(k)} = 0$, $k = 1, 2, \dots, p$ for all orders up to and including p . Assumptions $C(\eta)$ and $D(\zeta)$ are sometimes referred to as the row and column simplifying assumptions, respectively. The stage-order of a Runge–Kutta method is the largest value of q such that $B(q)$ and $C(q)$ are both satisfied. As its name implies, stage-order is related to the order of accuracy of the intermediate stage values of the U -vector, U_i , and typically equals the lowest order amongst all stages.

Closely related to assumption $B(p)$, $p^{(k)}$ is defined as

$$p^{(k)} = \mathbf{b}^T \mathbf{C}^{k-1} \mathbf{e} - \frac{1}{k} = \tau_1^k (k-1)!, \quad (11)$$

where $\mathbf{C} = \text{diag}(\mathbf{c})$, $\mathbf{e} = \{1, 1, \dots, 1\}$ and $\mathbf{C}\mathbf{e} = \mathbf{c}$. Writing $p_1^{(1,2,3,4)} = 0$ is meant to imply that $p_1^{(1)} = p_1^{(2)} = p_1^{(3)} = p_1^{(4)} = 0$. Powers of the vector \mathbf{c} should be interpreted as componentwise multiplication or, the Hadamard product. Hence $\mathbf{c}^3 = \mathbf{c} * \mathbf{c} * \mathbf{c} = \mathbf{C}^3 \mathbf{e}$ and $\mathbf{c}^0 = \mathbf{e}$. Powers of the \mathbf{A} -matrix are given by $\mathbf{A}^0 = \mathbf{I}$, $\mathbf{A}^1 = \mathbf{A}$, $\mathbf{A}^2 = \mathbf{AA}$ etcetera.

Similarly, $\mathbf{q}^{(k)}$ is closely related to assumption $C(\eta)$,

$$\mathbf{q}^{(k)} = \mathbf{AC}^{k-1} \mathbf{e} - \frac{1}{k} \mathbf{C}^k \mathbf{e}, \quad \mathbf{q}^{(1)} = \mathbf{A}\mathbf{e} - \mathbf{c}, \quad \text{Diag}(\mathbf{q}^{(k)}) = \mathbf{Q}^{(k)}, \quad (12)$$

where $\mathbf{q}^{(1)} = \mathbf{0}$ is simply the row-sum condition. Writing $q_{2,3,4,5,6}^{(k)} = 0$ is meant to imply that $q_2^{(k)} = q_3^{(k)} = q_4^{(k)} = q_5^{(k)} = q_6^{(k)} = 0$. Table 1 details the use of $\mathbf{q}^{(k)}$, and consequently $C(\eta)$, for order conditions up to order six. When one attempts to apply

Table 1

Order conditions expressed using $\mathbf{q}^{(k)}$ and $\mathbf{r}^{(k)}$ up to sixth-order for Runge–Kutta methods. Bushy tree order conditions, $\tau_1^{(l)}, l = 1, 2, \dots, 6$, are not included.

$C(\eta)$	$D(\zeta)$
$\tau_2^{(3)}$	$\mathbf{b}^T \mathbf{q}^{(2)} + \tau_1^{(3)}$
$\tau_2^{(4)}$	$\mathbf{b}^T \mathbf{Cq}^{(2)} + 3\tau_1^{(4)}$
$\tau_3^{(4)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{q}^{(3)} + \tau_1^{(4)}$
$\tau_4^{(4)}$	$\mathbf{b}^T \mathbf{Aq}^{(2)} + \tau_3^{(4)}$
$\tau_2^{(5)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{C}^2 \mathbf{q}^{(2)} + 6\tau_1^{(5)}$
$\tau_3^{(5)}$	$\frac{1}{2}\mathbf{b}^T (\mathbf{q}^{(2)} + \mathbf{C}^2) \mathbf{q}^{(2)} + 3\tau_1^{(5)}$
$\tau_4^{(5)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{Cq}^{(3)} + 4\tau_1^{(5)}$
$\tau_5^{(5)}$	$\frac{1}{6}\mathbf{b}^T \mathbf{q}^{(4)} + \tau_1^{(5)}$
$\tau_6^{(5)}$	$\mathbf{b}^T \mathbf{CAq}^{(2)} + \tau_4^{(5)}$
$\tau_7^{(5)}$	$\mathbf{b}^T \mathbf{ACq}^{(2)} + 3\tau_5^{(5)}$
$\tau_8^{(5)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{Aq}^{(3)} + \tau_5^{(5)}$
$\tau_9^{(5)}$	$\mathbf{b}^T \mathbf{AAq}^{(2)} + \tau_8^{(5)}$
$\tau_2^{(6)}$	$\frac{1}{6}\mathbf{b}^T \mathbf{C}^3 \mathbf{q}^{(2)} + 10\tau_1^{(6)}$
$\tau_3^{(6)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{C}(\mathbf{q}^{(2)} + \mathbf{C}^2) \mathbf{q}^{(2)} + 15\tau_1^{(6)}$
$\tau_4^{(6)}$	$\frac{1}{4}\mathbf{b}^T \mathbf{C}^2 \mathbf{q}^{(3)} + 10\tau_1^{(6)}$
$\tau_5^{(6)}$	$\frac{1}{2}\mathbf{b}^T (\mathbf{Q}^{(3)} + \frac{1}{3}\mathbf{C}^3) \mathbf{q}^{(2)} + \tau_4^{(6)}$
$\tau_6^{(6)}$	$\frac{1}{6}\mathbf{b}^T \mathbf{Cq}^{(4)} + 5\tau_1^{(6)}$
$\tau_7^{(6)}$	$\frac{1}{24}\mathbf{b}^T \mathbf{q}^{(5)} + \tau_1^{(6)}$
$\tau_8^{(6)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{C}^2 \mathbf{Aq}^{(2)} + \tau_4^{(6)}$
$\tau_9^{(6)}$	$\mathbf{b}^T \mathbf{Q}^{(2)} \mathbf{Aq}^{(2)} + \frac{1}{2}\mathbf{b}^T \mathbf{Q}^{(2)} \mathbf{q}^{(3)} + \frac{1}{6}\mathbf{b}^T \mathbf{C}^3 \mathbf{q}^{(2)} + \tau_8^{(6)}$
$\tau_{10}^{(6)}$	$\mathbf{b}^T \mathbf{CACq}^{(2)} + 3\tau_6^{(6)}$
$\tau_{11}^{(6)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{AC}^2 \mathbf{q}^{(2)} + 6\tau_7^{(6)}$
$\tau_{12}^{(6)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{A}(\mathbf{q}^{(2)} + \mathbf{C}^2) \mathbf{q}^{(2)} + 3\tau_7^{(6)}$
$\tau_{13}^{(6)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{CAq}^{(3)} + \tau_6^{(6)}$
$\tau_{14}^{(6)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{ACq}^{(3)} + 4\tau_7^{(6)}$
$\tau_{15}^{(6)}$	$\frac{1}{6}\mathbf{b}^T \mathbf{Aq}^{(4)} + \tau_7^{(6)}$
$\tau_{16}^{(6)}$	$\mathbf{b}^T \mathbf{CAAq}^{(2)} + \tau_{13}^{(6)}$
$\tau_{17}^{(6)}$	$\mathbf{b}^T \mathbf{ACAq}^{(2)} + \tau_{14}^{(6)}$
$\tau_{18}^{(6)}$	$\mathbf{b}^T \mathbf{AACq}^{(2)} + 3\tau_{15}^{(6)}$
$\tau_{19}^{(6)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{AAq}^{(3)} + \tau_{15}^{(6)}$
$\tau_{20}^{(6)}$	$\mathbf{b}^T \mathbf{AAAq}^{(2)} + \tau_{19}^{(6)}$

$\mathbf{q}^{(3)}$ to a stage-order two ESDIRK, it is found that $q_2^{(3)} \neq 0$. To denote this case, a truncated version of $C(3)^*$ is defined as $q_i^{(3)} = 0, i = 3, 4, \dots, s$.

Lastly, closely related to assumption $D(\zeta)$, one may further define

$$\mathbf{r}^{(k)} = \mathbf{b}^T \mathbf{C}^{k-1} \mathbf{A} - \frac{1}{k} \mathbf{b}^T (\mathbf{I} - \mathbf{C}^k). \quad (13)$$

This assumption principally applies to subquadrature and extended subquadrature order conditions as most of the nonlinear order conditions cannot be reduced by using the column simplifying assumption. Because methods in this paper will be stiffly-accurate, the application of any $D(\zeta)$ assumptions is unhelpful.

Lastly, simplifying assumptions $B(p)$, $C(\eta)$ and $D(\zeta)$ may be related through $p^{(k)}$, $\mathbf{q}^{(k)}$ and $\mathbf{r}^{(k)}$ as

$$\mathbf{r}^{(k)} \mathbf{C}^{l-1} = \mathbf{b}^T \mathbf{C}^{k-1} \mathbf{q}^{(l)} + \frac{(k+l)}{kl} p^{(k+l)} - \frac{1}{k} p^{(l)}. \quad (14)$$

For orders up to six, order conditions are expressed in Table 1 where $\mathbf{q}^{(k)}$ and $\mathbf{r}^{(k)}$ are both used. Table 2 shows the effect applying both $C(2)$ and $C(3)^*$ on all subquadrature, extended subquadrature and nonlinear order conditions up to sixth-order.

Table 2

Order conditions expressed using $\mathbf{q}^{(k)}$ up to sixth-order and their reduced from upon application of $C(2)$ or a truncated form of $C(3)^*$. Bushy tree order conditions, $\tau_1^{(l)}, l = 1, 2, \dots, 6$, are not included.

$C(\eta)$	$C(2)$	$C(3)^*$
$\tau_2^{(3)}$	$\mathbf{b}^T \mathbf{q}^{(2)} + \tau_1^{(3)}$	$\tau_1^{(3)}$
$\tau_2^{(4)}$	$\mathbf{b}^T \mathbf{Cq}^{(2)} + 3\tau_1^{(4)}$	$3\tau_1^{(4)}$
$\tau_3^{(4)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{q}^{(3)} + \tau_1^{(4)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{q}^{(3)} + \tau_1^{(4)}$
$\tau_4^{(4)}$	$\mathbf{b}^T \mathbf{Aq}^{(2)} + \tau_3^{(4)}$	$\tau_3^{(4)}$
$\tau_5^{(5)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{C}^2 \mathbf{q}^{(2)} + 6\tau_1^{(5)}$	$6\tau_1^{(5)}$
$\tau_5^{(5)}$	$\frac{1}{2}\mathbf{b}^T (\mathbf{q}^{(2)} + \mathbf{C}^2) \mathbf{q}^{(2)} + 3\tau_1^{(5)}$	$3\tau_1^{(5)}$
$\tau_4^{(5)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{Cq}^{(3)} + 4\tau_1^{(5)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{Cq}^{(3)} + 4\tau_1^{(5)}$
$\tau_5^{(5)}$	$\frac{1}{6}\mathbf{b}^T \mathbf{q}^{(4)} + \tau_1^{(5)}$	$\frac{1}{6}\mathbf{b}^T \mathbf{q}^{(4)} + \tau_1^{(5)}$
$\tau_6^{(5)}$	$\mathbf{b}^T \mathbf{CAq}^{(2)} + \tau_4^{(5)}$	$\tau_4^{(5)}$
$\tau_7^{(5)}$	$\mathbf{b}^T \mathbf{ACq}^{(2)} + 3\tau_5^{(5)}$	$3\tau_5^{(5)}$
$\tau_8^{(5)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{Aq}^{(3)} + \tau_5^{(5)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{Aq}^{(3)} + \tau_5^{(5)}$
$\tau_9^{(5)}$	$\mathbf{b}^T \mathbf{AAq}^{(2)} + \tau_8^{(5)}$	$\tau_8^{(5)}$
$\tau_{10}^{(6)}$	$\frac{1}{6}\mathbf{b}^T \mathbf{C}^3 \mathbf{q}^{(2)} + 10\tau_1^{(6)}$	$10\tau_1^{(6)}$
$\tau_{10}^{(6)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{C}(\mathbf{q}^{(2)} + \mathbf{C}^2) \mathbf{q}^{(2)} + 15\tau_1^{(6)}$	$15\tau_1^{(6)}$
$\tau_{11}^{(6)}$	$\frac{1}{4}\mathbf{b}^T \mathbf{C}^2 \mathbf{q}^{(3)} + 10\tau_1^{(6)}$	$\frac{1}{4}\mathbf{b}^T \mathbf{C}^2 \mathbf{q}^{(3)} + 10\tau_1^{(6)}$
$\tau_{12}^{(6)}$	$\frac{1}{2}\mathbf{b}^T (\mathbf{Q}^{(3)} + \frac{1}{2}\mathbf{C}^3) \mathbf{q}^{(2)} + \tau_4^{(6)}$	$\tau_4^{(6)}$
$\tau_{13}^{(6)}$	$\frac{1}{6}\mathbf{b}^T \mathbf{Cq}^{(4)} + 5\tau_1^{(6)}$	$\frac{1}{6}\mathbf{b}^T \mathbf{Cq}^{(4)} + 5\tau_1^{(6)}$
$\tau_{14}^{(6)}$	$\frac{1}{24}\mathbf{b}^T \mathbf{q}^{(5)} + \tau_1^{(6)}$	$\frac{1}{24}\mathbf{b}^T \mathbf{q}^{(5)} + \tau_1^{(6)}$
$\tau_{15}^{(6)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{C}^2 \mathbf{Aq}^{(2)} + \tau_4^{(6)}$	$\tau_4^{(6)}$
$\tau_{16}^{(6)}$	$\mathbf{b}^T \mathbf{Q}^{(2)} \mathbf{Aq}^{(2)} + \frac{1}{2}\mathbf{b}^T \mathbf{Q}^{(2)} \mathbf{q}^{(3)} + \frac{1}{6}\mathbf{b}^T \mathbf{C}^3 \mathbf{q}^{(2)} + \tau_8^{(6)}$	$\tau_8^{(6)}$
$\tau_{17}^{(6)}$	$\mathbf{b}^T \mathbf{CACq}^{(2)} + 3\tau_6^{(6)}$	$3\tau_6^{(6)}$
$\tau_{18}^{(6)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{AC}^2 \mathbf{q}^{(2)} + 6\tau_7^{(6)}$	$6\tau_7^{(6)}$
$\tau_{19}^{(6)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{A}(\mathbf{q}^{(2)} + \mathbf{C}^2) \mathbf{q}^{(2)} + 3\tau_7^{(6)}$	$3\tau_7^{(6)}$
$\tau_{20}^{(6)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{CAq}^{(3)} + \tau_6^{(6)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{CAq}^{(3)} + \tau_6^{(6)}$
$\tau_{21}^{(6)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{ACq}^{(3)} + 4\tau_7^{(6)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{ACq}^{(3)} + 4\tau_7^{(6)}$
$\tau_{22}^{(6)}$	$\frac{1}{6}\mathbf{b}^T \mathbf{Aq}^{(4)} + \tau_7^{(6)}$	$\frac{1}{6}\mathbf{b}^T \mathbf{Aq}^{(4)} + \tau_7^{(6)}$
$\tau_{23}^{(6)}$	$\mathbf{b}^T \mathbf{CAAq}^{(2)} + \tau_{13}^{(6)}$	$\tau_{13}^{(6)}$
$\tau_{24}^{(6)}$	$\mathbf{b}^T \mathbf{ACAq}^{(2)} + \tau_{14}^{(6)}$	$\tau_{14}^{(6)}$
$\tau_{25}^{(6)}$	$\mathbf{b}^T \mathbf{AACq}^{(2)} + 3\tau_{15}^{(6)}$	$3\tau_{15}^{(6)}$
$\tau_{26}^{(6)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{AAq}^{(3)} + \tau_{15}^{(6)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{AAq}^{(3)} + \tau_{15}^{(6)}$
$\tau_{27}^{(6)}$	$\mathbf{b}^T \mathbf{AAAq}^{(2)} + \tau_{19}^{(6)}$	$\frac{1}{2}\mathbf{b}^T \mathbf{AAAq}^{(2)} + \tau_{19}^{(6)}$

3.2. Error

Error in a p th-order Runge–Kutta scheme may be quantified in a general way by taking the L_2 principal error norm,

$$A^{(p+1)} = \|\boldsymbol{\tau}^{(p+1)}\|_2 = \sqrt{\sum_{j=1}^{a_{p+1}} \left(\tau_j^{(p+1)}\right)^2}. \quad (15)$$

For single-step, embedded schemes where $\hat{p} = p - 1$, additional definitions are useful such as

$$\hat{\boldsymbol{\tau}}_k^{(\hat{p})} = \frac{1}{\sigma} \sum_{i=1}^s \hat{b}_i \Phi_{i,k}^{(\hat{p})} - \frac{\alpha}{\hat{p}!}, \quad \hat{A}^{(\hat{p}+1)} = \|\hat{\boldsymbol{\tau}}^{(\hat{p}+1)}\|_2, \quad (16)$$

$$B^{(\hat{p}+2)} = \frac{\hat{A}^{(\hat{p}+2)}}{\hat{A}^{(\hat{p}+1)}}, \quad C^{(\hat{p}+2)} = \frac{\|\hat{\boldsymbol{\tau}}^{(\hat{p}+2)} - \boldsymbol{\tau}^{(\hat{p}+2)}\|_2}{\hat{A}^{(\hat{p}+1)}}, \quad E^{(\hat{p}+2)} = \frac{A^{(\hat{p}+2)}}{\hat{A}^{(\hat{p}+1)}}, \quad (17)$$

and $D = \text{Max}\{|a_{ij}|, |b_i|, |\hat{b}_i|, |c_i|\}$ where the superscript circumflex denotes the values with respect to the embedded method. The order of the method, p , refers to the global order-of-accuracy while the local order-of-accuracy is given by $p + 1$.

Table 3

Bounds on γ for A-stable SDIRKs and ESDIRKs from orders two to nine where p is the order of accuracy and s_I is the number of implicit stages.

s_I, p	A-stable
2, 2	$1/4 \leq \gamma \leq \infty$
2, 3	$\gamma = (3 + \sqrt{3})/6$
3, 3	$1/3 \leq \gamma \leq 1.068579021301628806418834$
3, 4	$\gamma = 1.068579021301628806418834$
4, 4	$0.3943375672974064411272872 \leq \gamma \leq 1.2805797612753054573024841$
4, 5	–
5, 5	$0.2465051931428202746001423 \leq \gamma \leq 0.3618033988749894848204587$ $0.4207825127659933063870173 \leq \gamma \leq 0.4732683912582953244555885$ $\gamma = 0.4732683912582953244555885$
5, 6	
6, 6	$0.2840646380117982930387010 \leq \gamma \leq 0.5409068780733081049137798$
6, 7	–
7, 7	–
7, 8	–
8, 8	$0.2170497430943030918315779 \leq \gamma \leq 0.2647142465800596850440755$
8, 9	–

4. Stability

4.1. Linear stability

Linear stability of DIRK-type methods applied to ODEs is studied based on the equation $U' = \lambda U$ by using the stability function

$$R(z) = \frac{P(z)}{Q(z)} = \frac{\sum_{i=1}^s p_i z^i}{\sum_{i=1}^s q_i z^i} = \frac{\text{Det} [\mathbf{I} - z\mathbf{A} + z\mathbf{e} \otimes \mathbf{b}^T]}{\text{Det} [\mathbf{I} - z\mathbf{A}]} = 1 + z\mathbf{b}^T [\mathbf{I} - z\mathbf{A}]^{-1} \mathbf{e}, \quad (18)$$

where \mathbf{I} is the identity matrix and $z = \lambda \Delta t$. Similar expressions may be written for the embedded and dense output methods by simply replacing \mathbf{b} in (18) with $\hat{\mathbf{b}}$ and $\mathbf{b}(\theta)$, respectively. A method is called A-stable and its stability function is called A-acceptable if $|R(z)| \leq 1$ for $\Re(z) \leq 0$. If, in addition to A-stability, $R(z)_{z \rightarrow -\infty} = 0$, then the method is called L-stable, and its stability function is called L-acceptable. L-acceptable stability functions have $\deg Q(z) > \deg P(z)$. For SDIRKs and ESDIRKs, $Q(z) = s_I$ so that $p_i = 0$ for $i \geq s_I$ is necessary for an L-acceptable stability function.

A Runge-Kutta method is imaginary axis or I-stable if $|R(z)| \leq 1$ for $\Re(z) = 0$. To test for I-stability, the E-polynomial is used,

$$E(y) = Q(+iy)Q(-iy) - P(+iy)P(-iy) = \sum_{j=0}^s E_{2j} y^{2j}, \quad (19)$$

where $i = \sqrt{-1}$. I-stability requires that $E(y) \geq 0$ for all values of real y , which requires $E(y)$ to have only imaginary roots. With the E-polynomial, one may establish a priori the bounds on γ that will result in either A-stable or L-stable methods provided that $p \approx s$. These bounds are given in Tables 3 and 4. To build these tables, methods were tested up to order 25 for both A- and L-stable methods. When $p \ll s$, enforcing $E(y) \geq 0$ can be difficult although γ can typically be reduced considerably relative to methods with $s - 1$ stages. Increasing the number of stages increases the number of right-hand-side evaluations as well as the required number of iterative solves but facilitates those iterative solves by reducing the condition number of the iteration matrix.

4.2. Nonlinear stability

Defining the symmetric algebraic stability matrix as

$$\mathbf{M} = \mathbf{B}^T \mathbf{A} + \mathbf{A}^T \mathbf{B} - \mathbf{b} \mathbf{b}^T, \quad M_{ij} = b_i a_{ij} + b_j a_{ji} - b_i b_j, \quad (20)$$

where $\mathbf{B} = \text{diag}(\mathbf{b})$ or $\mathbf{b} = \mathbf{B}\mathbf{e}$. Algebraic-stability of irreducible methods requires that $a_{ii}, b_i > 0$ and $\mathbf{M} \geq 0$. Consequently, SDIRKs may be algebraically stable but ESDIRKs may not. In [14], page 26, the term $\mathbf{B}\mathbf{B}^T$ should be replaced with $\mathbf{b}\mathbf{b}^T$.

Table 4

Bounds on γ for L-stable SDIRKs and ESDIRKs from orders two to eleven where p is the order of accuracy and s_l is the number of implicit stages.

s_l, p	L-stable
2, 2	$\gamma = (2 \pm \sqrt{2})/2$
3, 2	$0.1804253064293985641345831 \leq \gamma \leq 2.1856000973550400826291400$
3, 3	$\gamma = 0.43586652150845899941601945$
4, 3	$0.2236478009341764510696898 \leq \gamma \leq 0.5728160624821348554080014$
4, 4	$\gamma = 0.5728160624821348554080014$
5, 4	$0.2479946362127474551679910 \leq \gamma \leq 0.6760423932262813288723863$
5, 5	$\gamma = 0.2780538411364523249315862$
6, 5	$0.1839146536751751632321436 \leq \gamma \leq 0.3341423670680504359540301$
6, 6	$\gamma = 0.3341423670680504359540301$
7, 6	$0.2040834517158857633717906 \leq \gamma \leq 0.3788648944853283440258853$
7, 7	—
8, 7	$0.1566585993970439483924506 \leq \gamma \leq 0.2029348608433776737779349$ $0.2051941719494007117460614 \leq \gamma \leq 0.2343731596055835579475589$
8, 8	$\gamma = 0.2343731596055835579475589$
9, 8	$0.1708919625574635309332223 \leq \gamma \leq 0.2594205104814425547669495$
9, 9	—
10, 9	—
10, 10	—
11, 10	$0.1468989308591125260680428 \leq \gamma \leq 0.1657926100980560571096175$ $0.1937733662800920635754554 \leq \gamma \leq 0.1961524231108803003116274$
11, 11	—

4.3. Internal stability

Beyond traditional stepwise stability, it may be useful to control the stability associated with each stage in addition to each step. This is particularly true for large scaled eigenvalues, z , associated with stiff problems. To determine the vector of internal stabilities of Runge–Kutta methods, one evaluates

$$\begin{aligned} R_{\text{int}}(z) &= (\mathbf{I} - z\mathbf{A})^{-1}\mathbf{e} = \{R_{\text{int}}^{(1)}(z), R_{\text{int}}^{(2)}(z), \dots, R_{\text{int}}^{(s)}(z)\}^T, \\ &= \left\{ \frac{P_{\text{int}}^{(1)}(z)}{Q_{\text{int}}^{(1)}(z)}, \frac{P_{\text{int}}^{(2)}(z)}{Q_{\text{int}}^{(2)}(z)}, \dots, \frac{P_{\text{int}}^{(s)}(z)}{Q_{\text{int}}^{(s)}(z)} \right\}^T. \end{aligned} \quad (21)$$

The primary concern will be the value of $R_{\text{int}}^{(i)}(-\infty)$. One may also consider the E-polynomial, (19), at internal stages to determine stagewise I-stability by using

$$E_{\text{int}}^{(i)}(y) = Q_{\text{int}}^{(i)}(iy)Q_{\text{int}}^{(i)}(-iy) - P_{\text{int}}^{(i)}(iy)P_{\text{int}}^{(i)}(-iy), \quad (22)$$

where $E_{\text{int}}^{(i)}(y) \geq 0$ implies stagewise I-stability. One could also consider a stagewise analog to the algebraic-stability matrix for irreducible methods where

$$M_{jk}^{(i)} = a_{ij}a_{jk} + a_{ik}a_{kj} - a_{ij}a_{ik} \geq 0, \quad a_{ij} \geq 0, \quad i, j, k = 1, 2, \dots, s, \quad (23)$$

and $M_{jk}^{(i)}$ is the internal algebraic-stability matrix for stage i . Interestingly, the classic fully-implicit Runge–Kutta methods such as the Gauss, Radau (IA and IIA) and Lobatto (IIIA, IIIB and IIIC) methods are not internally algebraically stable on all stages. This suggests that internal algebraic stability is a rather severe requirement.

5. Implementation

For small systems of differential equations, one generally solves the implicit algebraic equations directly. In this case, using the definition $F_j = F(U_j, t^{[n]} + c_j \Delta t)$, one must solve

$$U_i = U^{[n]} + X_i + (\Delta t)\gamma F_i, \quad X_i = (\Delta t) \sum_{j=1}^{i-1} a_{ij} F_j, \quad 1 \leq i \leq s, \quad (24)$$

Table 5
Error controller coefficients.

Controller	α	β	γ	a	b
I = H ₀ 110	$\frac{1}{\hat{p}+1}$	0	0	0	0
H211	$\frac{1}{4\hat{p}}$	$\frac{-1}{4\hat{p}}$	0	$\frac{-1}{4}$	0
PC = H ₀ 220	$\frac{2}{\hat{p}}$	$\frac{1}{\hat{p}}$	0	1	0
PID	$\frac{1}{18\hat{p}}$	$\frac{-1}{9\hat{p}}$	$\frac{1}{18\hat{p}}$	0	0
H312	$\frac{1}{8\hat{p}}$	$\frac{-1}{4\hat{p}}$	$\frac{1}{8\hat{p}}$	$\frac{-3}{8}$	$\frac{-1}{8}$
PPIID	$\frac{6}{20\hat{p}}$	$\frac{-1}{20\hat{p}}$	$\frac{-5}{20\hat{p}}$	1	0
H321	$\frac{1}{3\hat{p}}$	$\frac{-1}{18\hat{p}}$	$\frac{-5}{18\hat{p}}$	$\frac{5}{6}$	$\frac{1}{6}$

where X_i is explicitly computed from existing data. Combining (24) with an appropriate starting guess, a modified Newton iteration provides U_i and F_i . This is accomplished by solving

$$(\mathbf{I} - \gamma(\Delta t)\mathbf{J})(U_{i,k+1} - U_{i,k}) = -(U_{i,k} - U^{[n]}) + X_i + \gamma(\Delta t)F_i, \quad (25)$$

where the subscript k denotes the value on the k th iteration, $(\partial F/\partial U) = \mathbf{J}$ is the Jacobian, \mathbf{I} is the identity matrix, $(\mathbf{I} - \gamma(\Delta t)\mathbf{J}) = \mathbf{N}$ is the (Newton) iteration matrix, and $(U_{i,k+1} - U_{i,k}) = \mathbf{d}_{i,k}$ is the displacement vector. The RHS of (25) is called the residual, $\mathbf{r}_{i,k}$ where $\mathbf{N}\mathbf{d}_{i,k} = \mathbf{r}_{i,k}$. Solving for the displacement vector, $U_{i,k+1}$ is then updated. Ideally, this procedure is repeated until some convergence criterion has been met. The importance of reducing the value of γ can be seen by looking at the term $\gamma(\Delta t)$ within the iteration matrix. Keeping the condition number of iteration matrix the same, a reduction in γ then permits an increase in the step size.

6. Step-size control

Local integration error for Runge–Kutta methods is usually controlled by first creating a local error estimate via an embedded method. This error estimate is then fed to an error controller which adjusts the time step in order to maintain some user-specified relative error tolerance ϵ . Step-size controllers are considered of the form

$$(\Delta t)^{[n+1]} = \kappa(\Delta t)^{[n]} \left\{ \frac{\epsilon}{||\delta^{[n+1]}||} \right\}^\alpha \left\{ \frac{||\delta^{[n]}||}{\epsilon} \right\}^\beta \left\{ \frac{\epsilon}{||\delta^{[n-1]}||} \right\}^\gamma \left\{ \frac{(\Delta t)^{[n]}}{(\Delta t)^{[n-1]}} \right\}^a \left\{ \frac{(\Delta t)^{[n-1]}}{(\Delta t)^{[n-2]}} \right\}^b, \quad (26)$$

for $p(\hat{p})$ -pairs ($\hat{p} = p - 1$). In (26), $\kappa \approx 0.95$ is the safety factor, $(\Delta t)^{[n-i]} = t^{[n-i]} - t^{[n-i-1]}$ is the step size, and $\delta^{[n+1]}$ is the vector of most recent local error estimates of the integration from $t^{[n-1]}$ to $t^{[n]}$ associated with the computation of $U^{[n+1]}$. Many controllers exist and Table 5 lists coefficients for various controllers.

7. Third-order methods

7.1. Four stages, $s_I = 3$

Stiffly-accurate, stage-order two ESDIRK methods in four stages appear as

0	0	0	0	
2γ	γ	γ	0	0
c_3	a_{31}	a_{32}	γ	0
1	b_1	b_2	b_3	γ
b_i	b_1	b_2	b_3	γ
\hat{b}_i	\hat{b}_1	\hat{b}_2	\hat{b}_3	\hat{b}_4

(27)

To enforce third-order overall accuracy for ODEs, a stage-order of two, and L-stability, one must satisfy the following six equations: $\tau_1^{(1,2,3)} = p_3 = q_{2,3}^{(2)} = 0$, leaving a one-parameter family of methods in c_3 . The general solution is

$$a_{32} = \frac{c_3(c_3 - 2\gamma)}{4\gamma}, \quad b_2 = \frac{-2 + 3c_3 + 6\gamma(1 - c_3)}{12\gamma(c_3 - 2\gamma)}, \quad b_3 = \frac{1 - 6\gamma + 6\gamma^2}{3c_3(c_3 - 2\gamma)}, \quad (28)$$

where $c_3 - 2\gamma \neq 0$ and $\gamma, c_3 \neq 0$. Demanding L-stability of the step requires that $\gamma = 0.43586652150845899941601945$. Selecting $c_3 = 3/5$, a method is found having $A^{(4)} = 0.03663$ and $R_{\text{int}}^{(3)}(-\infty) = -0.8057$. Stages 2 and 3 are both internally L-stable. Second-order, A-stable, error-control is accomplished by solving $\hat{\tau}_1^{(1,2)} = \hat{p}_4 = 0$, $\hat{R}(-\infty) = -\hat{p}_3/\gamma^3 = \gamma/2$. ESDIRK3(2)4L[2]SA was given in [14] and its properties are listed in Table 6.

Table 6
Third-order methods.

Name	ESDIRK 3(2)4L[2]SA	ESDIRK 3(2)5L[2]SA
s	4	5
p	3	3
γ	0.4359	$\frac{9}{40}$
$A^{(4)}$	0.03663	0.000777
$A^{(5)}$	0.07870	0.005199
$A^{(6)}$	0.1192	0.007633
$\hat{A}^{(3)}$	0.02552	0.002357
$\hat{A}^{(4)}$	0.07418	0.002437
$\{B^{(5)}, C^{(5)}, E^{(5)}\}$	{2.90, 1.64, 1.44}	{1.03, 1.21, 0.330}
D	1.000	1.000
$\{\lambda_{\text{Min}}^M, \lambda_{\text{Min}}^{\hat{M}}\}$	{−1.133, −2.335}	{−0.484, −0.340}
$\{b_{i,\text{Min}}, a_{ij,\text{Min}}\}$	{−0.595, −0.595}	{−0.347, −0.347}
$\lambda_{\text{Min}}^{M^{(2)}}$	−0.190	−0.051
$\lambda_{\text{Min}}^{M^{(3)}}$	−0.099	−0.075
$\lambda_{\text{Min}}^{M^{(4)}}$	−1.133	−0.051
$\lambda_{\text{Min}}^{M^{(5)}}$	−	−0.484
$\{R(-\infty), \hat{R}(-\infty)\}$	{0.0, 0.2179}	{0.0, 0.0}

7.2. Five stages, $s_I = 4$

To design a stiffly-accurate, stage-order two ESDIRK 3(2)-pair in five-stages, the following conditions are enforced

$$0 = \tau_1^{(1,2,3)} = q_{2,3,4,5}^{(1)} = q_{2,3,4}^{(2)} = p_4 = R_{\text{int}}^{(3,4)}(-\infty) = \hat{\tau}_1^{(1,2)} = \hat{p}_{4,5}, \quad \hat{\tau}_1^{(3)} = \frac{1}{600} \quad (29)$$

leaving γ and c_4 in main method. L-stable methods may be found for

$$0.2236478009341764510696898 \leq \gamma \leq 0.5728160624821348554080014. \quad (30)$$

Setting $\gamma = 9/40 = 0.225$, $c_3 = 9(2 + \sqrt{2})/40$ and $c_4 = 3/5$ yield a method having $A^{(4)} = 0.0007769$. This represents over an order-of-magnitude reduction in both $A^{(4)}$ and $A^{(5)}$ relative to ESDIRK3(2)4L[2]SA as well as a reduction in γ by nearly 50%. The method is L-stable on all internal stages and, hence, A-stable on stage 2 and L-stable on stages 3, 4 and 5. Further, it has an L-stable embedded method. This scheme, ESDIRK3(2)5L[2]SA, was given in [14] while its properties are listed in Table 6. Minimum eigenvalues of the algebraic stability matrix, λ_{Min}^M , and the internal algebraic stability matrix in stage i , $\lambda_{\text{Min}}^{M^{(i)}}$, are similar between ESDIRK3(2)4L[2]SA and ESDIRK3(2)5L[2]SA.

8. Fourth-order methods

8.1. Five stages, $s_I = 5$

Arguably, the most popular DIRK-type method is SDIRK4 [9]. It is an SDIRK in five-stages which is fourth-order but has only stage-order one. Key properties for this reference method are given in Table 7. From this table, SDIRK4 can be seen to be strongly nonlinearly unstable with $\lambda_{\text{Min}}^M = -112.1$.

8.2. Six stages, $s_I = 5$

Six-stage, stage-order two, stiffly-accurate ESDIRK methods

0	0	0	0	0	0	
2γ	γ	γ	0	0	0	0
c_3	a_{31}	a_{32}	γ	0	0	0
c_4	a_{41}	a_{42}	a_{43}	γ	0	0
c_5	a_{51}	a_{52}	a_{53}	a_{54}	γ	0
1	b_1	b_2	b_3	b_4	b_5	γ
b_i	\hat{b}_1	\hat{b}_2	\hat{b}_3	\hat{b}_4	\hat{b}_5	γ
\hat{b}_i	\hat{b}_1	\hat{b}_2	\hat{b}_3	\hat{b}_4	\hat{b}_5	\hat{b}_6

Table 7
Fourth-order methods.

Name	SDIRK4	ESDIRK 4(3)6L[2]SA ₁	ESDIRK 4(3)6L[2]SA ₂	ESDIRK 4(3)7L[2]SA
s	5	6	6	7
p	4	4	4	4
γ	$\frac{1}{4}$	$\frac{1}{4}$	0.248	$\frac{1}{8}$
$A^{(5)}$	0.002504	0.001830	0.001686	0.000260
$A^{(6)}$	0.004511	0.003467	0.002893	0.001177
$\hat{A}^{(4)}$	0.01247	0.003187	0.003187	0.000301
$\hat{A}^{(5)}$	0.01638	0.004077	0.004319	0.000977
$\{B^{(5)}, C^{(5)}, E^{(5)}\}$	{1.31, 1.43, 0.201}	{1.28, 1.15, 0.574}	{1.36, 1.15, 0.529}	{3.24, 3.07, 0.861}
D	7.813	1.585	1.504	1.00
$\{\lambda_{\text{Min}}^M, \lambda_{\text{Min}}^{\hat{M}}\}$	{−112.1, −100.45}	{−0.197, −0.184}	{−0.174, −0.148}	{−1.990, −1.361}
$\{b_{i,\text{Min}}, a_{ij,\text{Min}}\}$	{−7.083, −7.083}	{−0.108, −0.727}	{−0.098, −0.690}	{−0.557, −0.557}
$\lambda_{\text{Min}}^{M^{(2)}}$	+0.063	−0.063	−0.062	−0.016
$\lambda_{\text{Min}}^{M^{(3)}}$	−0.022	−0.036	−0.035	−0.009
$\lambda_{\text{Min}}^{M^{(4)}}$	−0.029	−0.063	−0.075	−0.314
$\lambda_{\text{Min}}^{M^{(5)}}$	−112.1	−3.360	−3.033	−1.225
$\lambda_{\text{Min}}^{M^{(6)}}$	−	−0.197	−0.174	−0.204
$\lambda_{\text{Min}}^{M^{(7)}}$	−	−	−	−1.990
$\{R(-\infty), \hat{R}(-\infty)\}$	{0.0, 3.3}	{0.0, 0.0}	{0.0, 0.0}	{0.0, 0.0}

provide $s(s+1)/2 - 4 = 17$ degrees of freedom (DoF), as shown, where $q_2^{(1,2)}$ and $\tau_1^{(1)}$ have already been applied. L-stable methods may be found for

$$0.2479946362127474551679910 \leq \gamma \leq 0.6760423932262813288723863. \quad (32)$$

Fourth-order methods are obtained by solving

$$0 = \tau_1^{(1,2,3,4)} = q_{2,3,4,5}^{(1)} = q_{2,3,4,5}^{(2)} = \tau_3^{(4)} = p_5 = R_{\text{int}}^{(3,4,5)}(-\infty), \quad (33)$$

$$0 = \hat{\tau}_1^{(1,2,3)} = \hat{p}_{5,6}, \quad \hat{\tau}_1^{(4)} = \frac{1}{1000}, \quad (34)$$

for the main method and the embedded methods. Selecting $c_3 = \gamma(2 - \sqrt{2})$, three DoF remain, say c_4 , c_5 and γ . Two solutions are given by

- ESDIRK4(3)6L[2]SA₁: $\gamma = \frac{1}{4}$, $c_4 = \frac{5}{8}$, $c_5 = \frac{26}{25}$,
- ESDIRK4(3)6L[2]SA₂: $\gamma = \frac{248}{1000}$, $c_4 = \frac{1043}{1706}$, $c_5 = \frac{1361}{1300}$.

Both methods are strictly internally A-stable on stage two but internally L-stable on stages three through six. Also, both embedded methods are L-stable. The first method, ESDIRK4(3)6L[2]SA₁, is given in [14] while the second method, ESDIRK4(3)6L[2]SA₂, is a slight perturbation on ESDIRK4(3)6L[2]SA₁. Notice that one of abscissae, c_5 , is just outside the integration step. This could be problematic when integrating across temporal discontinuities. Properties for both of these methods are given in Table 7 and the coefficients to the method are given in Table 9 in Appendix A.

8.3. Seven stages, $s_I = 6$

Moving to seven stages, fourth-order methods are obtained by solving [15]

$$0 = \tau_1^{(1,2,3,4,5)} = q_{2,3,4,5,6}^{(1)} = q_{2,3,4,5,6}^{(2)} = \tau_3^{(4)} = \tau_{4,5}^{(5)} = p_6 = R_{\text{int}}^{(3,4,5,6)}(-\infty), \quad (35)$$

$$0 = \hat{\tau}_1^{(1,2,3)} = \hat{p}_{6,7} = \frac{\partial \hat{A}^{(4)}}{\partial \hat{b}_6}, \quad (36)$$

for the main method and for the embedded method. Four DoF remain, e.g. γ , c_4 , c_5 and c_6 in the main method while one remains in the embedded method. The dominant motivation for creating a seven-stage, fourth-order ESDIRKs is that γ can be driven to quite small values while the leading order error can be driven quite low. Setting $\gamma = \frac{1}{8}$, $\hat{b}_6 = \frac{19}{140}$ and

$$c_2 = \frac{1}{4}, \quad c_3 = (2 - \sqrt{2})/8, \quad c_4 = \frac{1}{2}, \quad c_5 = \frac{395}{567}, \quad c_6 = \frac{89}{126}, \quad c_7 = 1, \quad (37)$$

yields a method, ESDIRK4(3)7L[2]SA, which is L-stable, internally L-stable on stages three through seven and strictly A-stable on the second stage. Leading-order error has reduced six-fold from the similar six-stage method and a reduction in γ by nearly 50% relative to ESDIRK4(3)6L[2]SA₂. Key properties for ESDIRK4(3)7L[2]SA are given in Table 7. Coefficients to the method are given in Table 10 in Appendix A. In comparison to ARK4(3)6L[2]SA₁, ESDIRK4(3)6L[2]SA₂ has half the leading-order error. In comparison to ARK4(3)7L[2]SA₁, ESDIRK4(3)7L[2]SA₁ has 16% of the leading order error.

9. Fifth-order methods

9.1. Seven stages, $s_I = 6$

Seven-stage, stiffly-accurate ESDIRK methods take the form

$$\begin{array}{c|ccccccc} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 2\gamma & \gamma & \gamma & 0 & 0 & 0 & 0 & 0 \\ c_3 & a_{31} & a_{32} & \gamma & 0 & 0 & 0 & 0 \\ c_4 & a_{41} & a_{42} & a_{43} & \gamma & 0 & 0 & 0 \\ c_5 & a_{51} & a_{52} & a_{53} & a_{54} & \gamma & 0 & 0 \\ c_6 & a_{61} & a_{62} & a_{63} & a_{64} & a_{65} & \gamma & 0 \\ 1 & b_1 & b_2 & b_3 & b_4 & b_5 & b_6 & \gamma \\ \hline b_i & b_1 & b_2 & b_3 & b_4 & b_5 & b_6 & \gamma \\ \hat{b}_i & \hat{b}_1 & \hat{b}_2 & \hat{b}_3 & \hat{b}_4 & \hat{b}_5 & \hat{b}_6 & \hat{b}_7, \end{array} \quad (38)$$

provide 24 DoF, as shown, where $q_2^{(1,2)}$ and $\tau_1^{(1)}$ have already been applied. Row simplifying assumption C(2) ensures a stage-order two method. L-stable methods may be found for a range γ values

$$0.1839146536751751632321436 \leq \gamma \leq 0.3341423670680504359540301. \quad (39)$$

Fifth-order methods are obtained by solving

$$0 = \tau_1^{(1,2,3,4,5)} = q_{2,3,\dots,6}^{(1)} = q_{2,3,4,5,6}^{(2)} = \tau_3^{(4)} = \tau_{4,5,8}^{(5)} = p_6 = R_{\text{int}}^{(3,4,5,6)}(-\infty), \quad (40)$$

$$0 = \hat{\tau}_1^{(1,2,3,4)} = \hat{\tau}_3^{(4)} = \hat{p}_7, \quad (41)$$

for the main method and for the embedded methods. This leaves three remaining DoF in the main method e.g., c_4 , c_6 and γ and \hat{b}_7 in the embedded method. As γ is probably best chosen as small as possible, $\gamma = 0.184 = 23/125$ is selected. The coefficient \hat{b}_7 has been used to enforce $\hat{p}_6/\gamma^6 = \hat{R}(-\infty)$. Two solutions to these conditions are shown in Table 8 and given by

- ESDIRK5(4)7L[2]SA₁: $c_3 = (2 - \sqrt{2})\gamma$, $c_4 = \frac{13}{25}$, $c_5 = \frac{5906118540659}{9042400211275}$, $c_6 = \frac{26}{25}$,
- ESDIRK5(4)7L[2]SA₂: $c_3 = (2 + \sqrt{2})\gamma$, $c_4 = \frac{49}{353}$, $c_5 = \frac{3706679970760}{5295570149437}$, $c_6 = \frac{347}{382}$.

The first method, ESDIRK5(4)7L[2]SA₁, was given in [14], as ESDIRK5(4)7L[2]SA. While ESDIRK5(4)7L[2]SA₁ is not I-stable and strongly nonlinearly unstable on the sixth-stage, ESDIRK5(4)7L[2]SA₂ is I-stable and only mildly nonlinearly unstable on the same stage. Table 11 of Appendix A contains the coefficients to ESDIRK5(4)7L[2]SA₂.

9.2. Eight stages, $s_I = 7$

If another stage is added, fifth-order methods are obtained by solving

$$0 = \tau_1^{(1,2,3,4,5,6)} = q_{2,3,\dots,7}^{(1)} = q_{2,3,\dots,7}^{(2)} = \tau_3^{(4)} = \tau_{4,5,8}^{(5)} = \tau_{6,7}^{(6)} = p_7 = R_{\text{int}}^{(3,4,5,6,7)}(-\infty), \quad (42)$$

$$0 = \hat{\tau}_1^{(1,2,3,4)} = \hat{\tau}_3^{(4)} = \hat{p}_{7,8}, \quad (43)$$

for the main method and for the embedded method. As seen from Table 1, enforcement of $\tau_{1,6,7}^{(6)} = 0$ forces $\tau_{2,3,10,11,12}^{(6)} = 0$. Four abscissae, $c_{4,5,6,7}$, and γ remain to be specified. Setting

$$\gamma = \frac{1}{7}, \quad c_2 = \frac{2}{7}, \quad c_3 = \frac{2 + \sqrt{2}}{7}, \quad c_4 = \frac{150}{203}, \quad c_5 = \frac{27}{46}, \quad c_6 = \frac{473}{532}, \quad c_7 = \frac{30}{83}, \quad (44)$$

and $\hat{b}_8 = 36/233$, ESDIRK5(4)8L[2]SA is found. Its properties are shown in Table 8.

Though ESDIRKs and fully implicit Runge–Kutta methods are quite different, the latter serve as reference methods. In particular, the Radau IIA family of methods are both L- and algebraically-stable. The three-stage Radau IIA method is fifth-order

Table 8

Fifth- and sixth-order methods.

Name	ESDIRK 5(4)7L[2]SA ₁	ESDIRK 5(4)7L[2]SA ₂	ESDIRK 5(4)8L[2]SA	ESDIRK 6(5)9L[2]SA
s	7	7	8	9
p	5	5	5	6
γ	$\frac{23}{125}$	$\frac{23}{125}$	$\frac{1}{7}$	$\frac{2}{9}$
$A^{(6)}$	0.001846	0.001272	0.0004459	—
$A^{(7)}$	0.003154	0.002184	0.0007294	0.0005388
$\hat{A}^{(5)}$	0.002171	0.002047	0.0003205	—
$\hat{A}^{(6)}$	0.001501	0.001882	0.0006473	0.003797
$\hat{A}^{(7)}$	0.000983	0.001215	0.0009153	0.0005831
$\{B^{(p+1)}, C^{(p+1)}, E^{(p+1)}\}$	{0.692, 1.31, 0.850}	{0.920, 1.26, 0.621}	{2.02, 1.77, 1.39}	{1.54, 1.60, 0.142}
D	8.971	1.634	1.000	0.9883
$\{\lambda_{\text{Min}}^{\mathbf{M}}, \lambda_{\text{Min}}^{\hat{\mathbf{M}}}\}$	{−0.800, −0.764}	{−0.405, −0.445}	{−1.256, −1.443}	{−1.254, −7.867}
$\{b_{i,\text{Min}}, a_{ij,\text{Min}}\}$	{−0.076, −5.036}	{−0.162, −0.539}	{−0.449, −0.847}	{−0.485, −0.922}
$\lambda_{\text{Min}}^{\mathbf{M}^{(2)}}$	−0.034	−0.034	−0.020	−0.049
$\lambda_{\text{Min}}^{\mathbf{M}^{(3)}}$	−0.019	−0.050	−0.030	−0.028
$\lambda_{\text{Min}}^{\mathbf{M}^{(4)}}$	−0.012	−0.012	−0.033	−0.223
$\lambda_{\text{Min}}^{\mathbf{M}^{(5)}}$	−4.290	−1.443	−0.033	−1.696
$\lambda_{\text{Min}}^{\mathbf{M}^{(6)}}$	−131.4	−2.519	−0.711	−0.331
$\lambda_{\text{Min}}^{\mathbf{M}^{(7)}}$	−0.800	−0.405	−1.259	−0.396
$\lambda_{\text{Min}}^{\mathbf{M}^{(8)}}$	—	—	−1.256	−2.866
$\lambda_{\text{Min}}^{\mathbf{M}^{(9)}}$	—	—	—	−1.254
$E_{\text{int}}^{(i)}(y) < 0$ $\{R(-\infty), \hat{R}(-\infty)\}$	$i = 6$ $\{0.0, 0.35\}$	— $\{0.0, -0.25\}$	— $\{0.0, 0.0\}$	$i = 6, 7$ $\{0.0, 0.1\}$

accurate with a stage-order of three with leading-order errors of $A^{(6)} = 0.0009895$ and $A^{(7)} = 0.001714$. The minimum values of the internal algebraic stability matrix are $\{-0.05566, -0.01750, 0.0\}$ for stages $\{1, 2, 3\}$ so that only stage 3 is internally algebraically stable where the coefficients of \mathbf{A} are all positive. All stages satisfy $R_{\text{int}}^{(i)}(-\infty) = 0$ but only stages 2 and 3 are internally I-stable. Coefficients to ESDIRK5(4)8L[2]SA are given in Table 12 of Appendix A. In comparison to ARK5(4)8L[2]SA₂, ESDIRK5(4)8L[2]SA has 60% of the leading order error and γ has been reduced 35%.

10. Sixth-order methods

To design a nine-stage ($s = 9$), stiffly-accurate, 6(5)-pair using only $C(2)$, $c_1 = 0$, $c_9 = 1$, and $a_{9j} = b_j$ are chosen and the main method is obtained with the satisfaction of the following equations,

$$0 = \tau_1^{(1,2,3,4,5,6)} = q_{2,3,\dots,8}^{(1)} = q_{2,3,\dots,8}^{(2)} = \tau_3^{(4)} = \tau_5^{(5)} = \tau_{4,6,7,13,14,15,19}^{(6)}, \quad (45)$$

$$= p_8 = R_{\text{int}}^{(3,4,5,6,7,8)}(-\infty). \quad (46)$$

Minimally, the embedded method must satisfy

$$0 = \hat{\tau}_1^{(1,2,3,4,5)} = \hat{\tau}_3^{(4)} = \hat{\tau}_{4,5,8}^{(5)} = \hat{p}_9. \quad (47)$$

$$1 \leq ||\hat{p}_8/q_8||. \quad (48)$$

Instead, a less arduous path towards the creation of a stiffly-accurate, 6(5)-pair in nine-stages involves applying a truncated version of C(3). Again, with $c_1 = 0$, $c_9 = 1$, $a_{9j} = b_j$ are set and the following equations are solved,

$$0 = \tau_1^{(1,2,3,4,5,6)} = q_{2,3,\dots,8}^{(1)} = q_{2,3,\dots,8}^{(2)} = q_{3,4,\dots,8}^{(3)} = \tau_5^{(5)} = \tau_{6,7,15}^{(6)}, \quad (49)$$

$$= \mathbf{b}^T \mathbf{e}_2 = \mathbf{b}^T \mathbf{A} \mathbf{e}_2 = \mathbf{b}^T \mathbf{C} \mathbf{A} \mathbf{e}_2 = \mathbf{b}^T \mathbf{A} \mathbf{A} \mathbf{e}_2 = p_8 = R_{\text{int}}^{(4,5,6,7,8)}(-\infty), \quad (50)$$

where $\mathbf{e}_2 = \{0, 1, 0, 0, \dots, 0\}$ is a vector of length s . As a consequence of $q_2^{(2)} = 0$, $c_2 = 2\gamma$, and since $q_3^{(3)} = 0$, $c_3 = (3 \pm \sqrt{3})\gamma$. The stiffly-accurate assumption along with $p_8 = 0$ ensures that $R(-\infty) = 0$. Expressions for the internal stability function and E-polynomial on stage three are found to be

$$R_{\text{int}}^{(3)}(-\infty) = \frac{c_3^2 - 4c_3\gamma + 2\gamma^2}{2\gamma^2}, \quad E_{\text{int}}^{(3)} = \frac{c_3}{4} (4\gamma - c_3)(c_3 - 2\gamma)^2 y^4. \quad (51)$$

To enforce $|R_{\text{int}}^{(3)}(-\infty)| < 1$, $c_3 = (3 - \sqrt{3})\gamma$. In this case,

$$R_{\text{int}}^{(3)}(-\infty) = (1 - \sqrt{3}), \quad E_{\text{int}}^{(3)} = (-3 + 2\sqrt{3})\gamma^4 y^4 > 0. \quad (52)$$

Hence, stage-three is I-stable and A-stable. Enforcing $R_{\text{int}}^{(4)}(-\infty) = 0$ where

$$R_{\text{int}}^{(4)}(-\infty) = \frac{-c_4^3 + 9c_4^2\gamma - 18c_4\gamma^2 + 6\gamma^3}{6\gamma^3}, \quad (53)$$

gives $c_4 = \left\{ 3 + \sqrt{3}\sin(\theta) \pm 3\cos(\theta), 3 - 2\sqrt{3}\sin(\theta) \right\} \gamma$ where $\theta = \cot^{-1}(\sqrt{2})/3$. Only one of these three solutions, $c_4 = [3 - 2\sqrt{3}\sin(\theta)]\gamma \approx \frac{803}{350}\gamma$, permits an I-stable fourth-stage where

$$\begin{aligned} E_{\text{int}}^{(4)} &= (\alpha_4 + \beta_4 y^2) y^4, \\ \alpha_4 &= \frac{c_4}{12} (c_4^3 - 12c_4^2\gamma + 36c_4\gamma^2 - 24\gamma^3), \\ \beta_4 &= \frac{c_4}{36} (6\gamma - c_4)(c_4 - 3\gamma)(c_4^3 - 9c_4^2\gamma + 18c_4\gamma^2 - 12\gamma^3), \end{aligned} \quad (54)$$

and $E_{\text{int}}^{(4)} \geq 0$ for $\alpha_4, \beta_4 \geq 0$. For stage five, enforcing $R_{\text{int}}^{(5)}(-\infty) = 0$ allows the internal E-polynomial for stage five to be written as

$$\begin{aligned} E_{\text{int}}^{(5)} &= (\alpha_5 + \beta_5 y^2 + \gamma^8 y^4) y^4, \\ \alpha_5 &= \frac{1}{12} (c_5^4 - 16c_5^3\gamma + 72c_5^2\gamma^2 - 96c_5\gamma^3 + 24\gamma^4), \\ \beta_5 &= \frac{1}{36} (c_5^3 - 12c_5^2\gamma + 36c_5\gamma^2 - 36\gamma^3)(c_5^3 - 12c_5^2\gamma + 36c_5\gamma^2 - 12\gamma^3). \end{aligned} \quad (55)$$

It may be determined that stage-five can be I-stable and L-stable if $E_{\text{int}}^{(5)} \geq 0$ for $y \geq 0$. This occurs if

$$1.745761101158346575686816712518 \leq \frac{c_5}{\gamma} \leq 4.471316041664625579088042286147. \quad (56)$$

There are five remaining DoF, say, γ and three abscissae from $c_{5,6,7,8}$. These DoF must be chosen so as to ensure I-stability of stages five through eight and the step. Satisfying those requirements, minimization of leading order error and nonlinear instability can be sought.

For the embedded method, there are more requirements than there are embedded scheme weights since, minimally, the embedded method must satisfy

$$0 = \hat{\tau}_1^{(1,2,3,4,5)} = \hat{b}_2 = \hat{\tau}_5^{(5)} = \hat{\mathbf{b}}^T \mathbf{A} \mathbf{e}_2 = \hat{p}_9. \quad (57)$$

Since these conditions also apply to the main method, care must be taken to keep the main and embedded methods distinct. Further, $|\hat{R}(-\infty)| = |\hat{p}_8/q_8| < 1$ or, $|\hat{p}_8| < \gamma^8$, to enable an A-stable embedded method. Together, the complete method has four DoF, say, γ , one \hat{b} and two abscissae from $c_{5,6,7,8}$.

To solve for the coefficients to the method, $\gamma = 2/9$ was selected and \hat{b}_7 was chosen as a proxy for $\hat{R}(-\infty)$. Lastly, c_5 as chosen to be with the range for I-stability on stage five and c_8 was selected to be a number between 0 and 1. With these four coefficients selected, Two equations remained in two coefficients, c_6 and c_7 . We were unable to find a good method that was I-stable on all stages. However, good methods were found with mild linear instability on stages six and seven. On each of the unstable stages, a small sliver of instability on the complex left-half-plane exists, immediately adjacent to the imaginary axis, within $-0.025 < \Re(z) = 0$ and $-3 < \Im(z) < 3$. The maximum magnitude of the internal linear stability functions for these two stages are

$$\begin{aligned} \text{Max}(R_{\text{int}}^{(6)}) &= 1.00429, z = \{-0.00000036883, \pm 1.8198\}, \\ \text{Max}(R_{\text{int}}^{(7)}) &= 1.00146, z = \{-0.00000088009, \pm 1.4273\}. \end{aligned}$$

For ESDIRK6(5)9L[2]SA, with

$$c_2 = 2\gamma, \quad c_3 = (3 - \sqrt{3})\gamma, \quad c_4 = [3 - 2\sqrt{3}\sin(\theta)]\gamma, \quad \theta = \cot^{-1}(\sqrt{2})/3, \quad (58)$$

the choices

$$\gamma = \frac{2}{9}, \quad c_5 = \frac{183}{200}, \quad c_8 = \frac{97}{100}, \quad \hat{R}(-\infty) = \frac{\hat{p}_8}{\gamma^8} = \frac{1}{10}, \quad (59)$$

give a method with

$$c_6 = \frac{62409086037595}{296036819031271}, \quad c_7 = \frac{81796628710131}{911762868125288}. \quad (60)$$

Its properties are shown in Table 8. Stages 4, 5, 8 and 9 are L-stable while stage 2 is strictly A-stable and stage 3 is strongly A-stable. Stages 6 and 7 are not I-stable but satisfy $R_{\text{int}}^{(6,7)}(-\infty) = 0$. Coefficients to this method are given in Appendix A, Table 13.

11. Test problems

Testing of schemes is conducted on two singular perturbation problems, van der Pol's (vdP) equation and Kaps' problem.

11.1. Van der Pol's equation

In 1926, Balthasar van der Pol wrote the equations describing a linear oscillating circuit. The van der Pol's (vdP) equation may be written as

$$\dot{y}_1 = y_2, \quad \dot{y}_2 = \varepsilon^{-1} \left[(1 - y_1^2) y_2 - y_1 \right] \quad (61)$$

It may be seen that for $\varepsilon \rightarrow 0$, the second equations become an algebraic relation, $y_2 = y_1 / (1 - y_1^2)$ (This corrects an error in [15]). This now constitutes an index-1 differential algebraic system where y_2 has emerged as the algebraic variable as $\varepsilon \rightarrow 0$. Unperturbed initials are given by [9]

$$y_1(0) = 2 \quad y_2(0) = -\frac{2}{3} + \frac{10}{81}\varepsilon - \frac{292}{2187}\varepsilon^2 + \frac{15266}{59049}\varepsilon^3 + \mathcal{O}(\varepsilon^4). \quad (62)$$

11.2. Kaps' problem

Dekker and Verwer [7] investigate a nonlinear problem (experiment 7.5.2) originally given by Peter Kaps,

$$\dot{y}_1 = -(\varepsilon^{-1} + 2)y_1 + \varepsilon^{-1}y_2^2, \quad \dot{y}_2 = y_1 - y_2 - y_2^2, \quad (63)$$

where $0 \leq t \leq 1$ and whose exact solution is $y_1 = y_2^2$, $y_2 = \exp(-t)$. As $\varepsilon \rightarrow 0$, the equations become

$$y_1 = y_2^2, \quad \dot{y}_2 = y_1 - y_2 - y_2^2, \quad (64)$$

which is an index-1 DAE. Equilibrium (unperturbed) initial conditions are given by $y_1(0) = y_2(0) = 1$.

For the stiffly accurate methods considered here, it may be anticipated that the observed orders of accuracy for the differential and algebraic variables while integrating these test problems are $h^p + \varepsilon h^{q+1}$ for the differential variable and $h^p + \varepsilon h^q$ for the algebraic variable where $h = (\Delta t)$ [9].

12. Discussion

The goal of the limited testing discussed below is to establish whether methods are appropriate for general purpose settings and are free of substantial shortcomings based on their convergence behavior on the two singular-perturbation test problems. Ten proposed and existing (E)SDIRK-type methods of third-, fourth-, fifth- and sixth-order accuracy are considered:

- Third-order: ESDIRK3(2)4L[2]SA and ESDIRK3(2)5L[2]SA,
- Fourth-order: SDIRK4, ESDIRK4(3)6L[2]SA₁ ESDIRK4(3)6L[2]SA₂ and ESDIRK4(3)7L[2]SA. Since ESDIRK4(3)6L[2]SA₁ and ESDIRK4(3)6L[2]SA₂ generate virtually identical test results, only the former will be shown.
- Fifth-order: ESDIRK5(4)7L[2]SA₁, ESDIRK5(4)7L[2]SA₂ and ESDIRK5(4)8L[2]SA
- Sixth-order: ESDIRK6(5)9L[2]SA

The two third-order methods, ESDIRK3(2)4L[2]SA and ESDIRK3(2)5L[2]SA, were previously tested [14]. ESDIRK3(2)5L[2]SA was found to be a better method. The methods ESDIRK4(3)6L[2]SA₁ and ESDIRK5(4)7L[2]SA₁ were also the subject of previous testing but are included here in the capacity of reference methods. SDIRK4 is included as it has, historically, been one of the most popular DIRK methods.

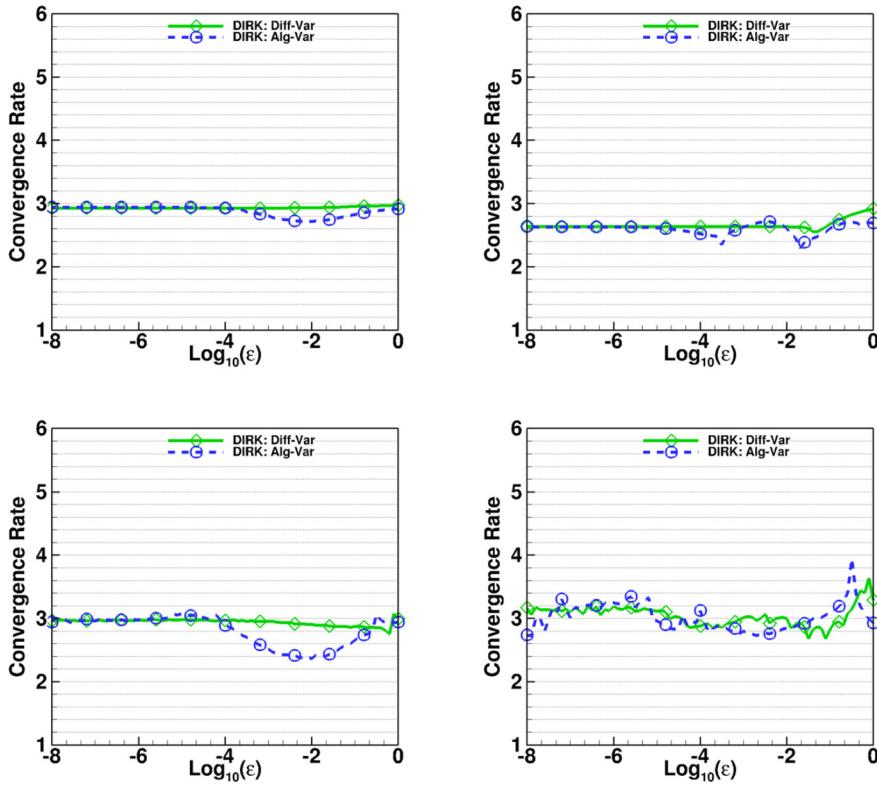


Fig. 1. A convergence rate comparison between ESDIRK3(2)4L[2]SA (left) and ESDIRK3(2)5L[2]SA (right) on Kaps' problem (top) and van der Pol's equation (bottom).

12.1. Convergence

All convergence tests were done using either the vdp or the Kaps' singular perturbation test problems, with the same initial conditions for each test problem, and integrating them over the same time interval. Step-sizes were held constant. Singular perturbation problems are very useful for detecting shortcomings of methods but may not be particularly illuminating at distinguishing which well-designed method is more accurate in practical settings. They also fail to exercise many important implementation matters. At orders three through six, comparisons between methods of similar orders will be made. These comparisons will consist of convergence-rate comparisons at various levels of the stiffness parameter, ϵ . Also, error versus work plots will be shown for the differential variable and emerging algebraic variables at selected values of ϵ . Jaggedness seen in many of the error-versus-work curves is likely driven by several factors including machine precision limitations. Convergence rates were inferred from curves using least-square fitting. Hence, many convergence rates inherit the nonsmooth character of their parents.

12.1.1. Third-order methods

Fig. 1 shows differences in convergence rates between ESDIRK3(2)4L[2]SA₂ and ESDIRK3(2)5L[2]SA₁ on both Kaps' problem and van der Pol's equation. Both methods have unremarkable behavior on both equations with convergence rates of essentially three for both the differential and emerging algebraic variables. For Kaps' problem, an error versus work comparison may be made between these same two methods at different values of the stiffness parameter. These are shown in Fig. 2 for the differential variable and algebraic variables. The five-stage method appears to be the more accurate method for a given level of work. Changes in the stiffness parameter appear to have minimal influence on the error versus work behavior of the methods. Similarly, for van der Pol's equation, Fig. 3 shows these same quantities. As opposed to the results using Kaps' problem, the influence of the stiffness parameter is more apparent on van der Pol's equation. Further, the five-stage method is significantly more sensitive to the choice of ϵ than is the four-stage method. In spite of this, it still appears to be the more efficient method. Given the much smaller value of γ with the five-stage method, it is also likely to show better iterative convergence than the four-stage method.

12.1.2. Fourth-order methods

Four of the eight methods tested in this paper are L-stable, stiffly-accurate, fourth-order methods. Of the four, ESDIRK4(3)6L[2]SA₂ was designed to eke out the last remaining performance from ESDIRK4(3)6L[2]SA₁. Principally, this

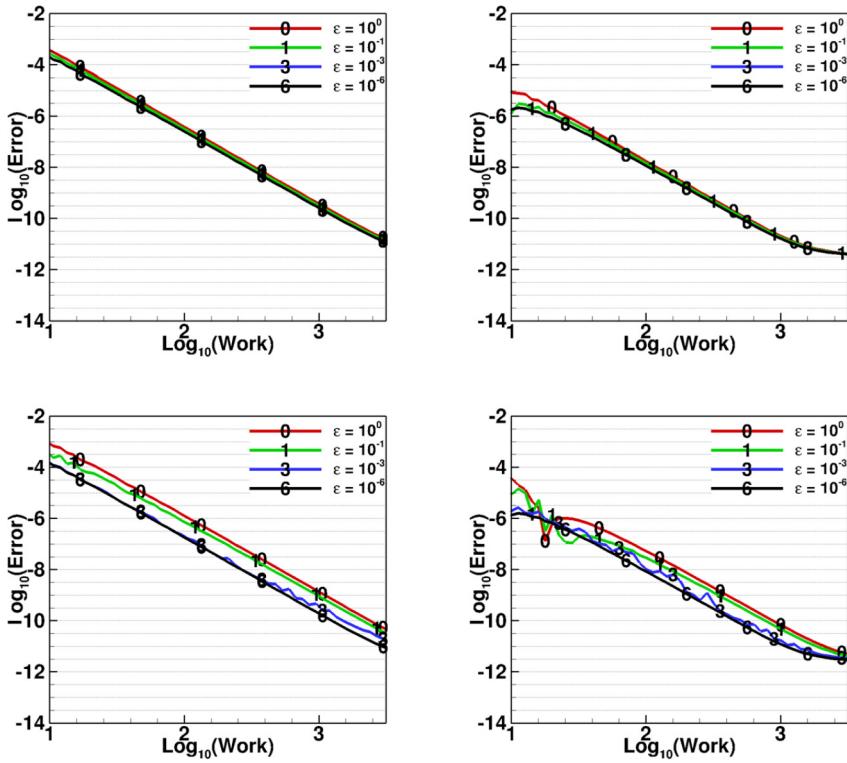


Fig. 2. An error versus work comparison between ESDIRK3(2)4L[2]SA (left) and ESDIRK3(2)5L[2]SA (right) for the differential variable (top) and algebraic variable (bottom) on Kaps' problem.

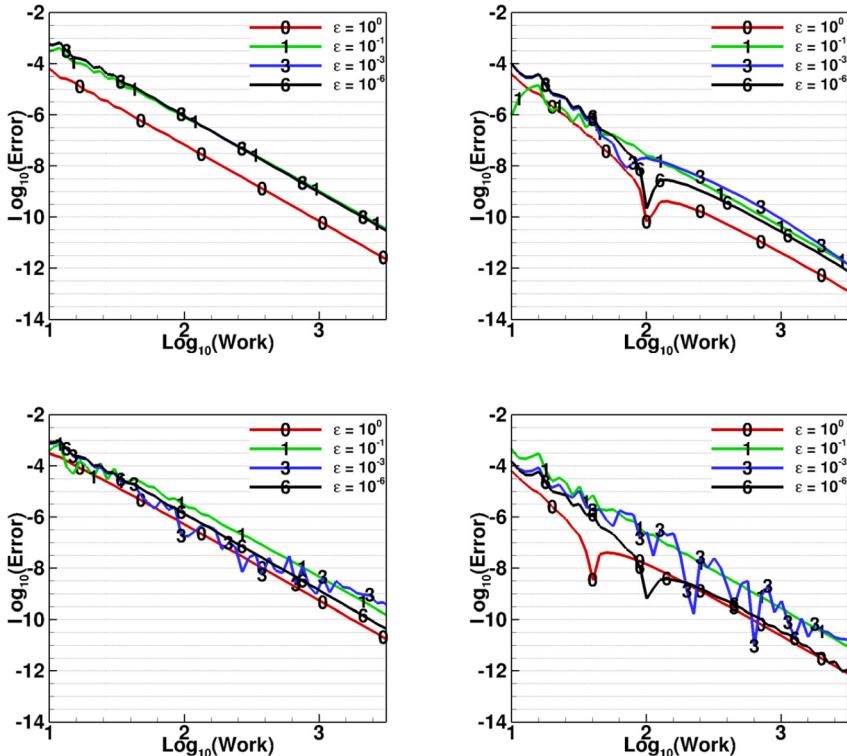


Fig. 3. An error versus work comparison between ESDIRK3(2)4L[2]SA (left) and ESDIRK3(2)5L[2]SA (right) for the differential variable (top) and algebraic variable (bottom) on van der Pol's equation.

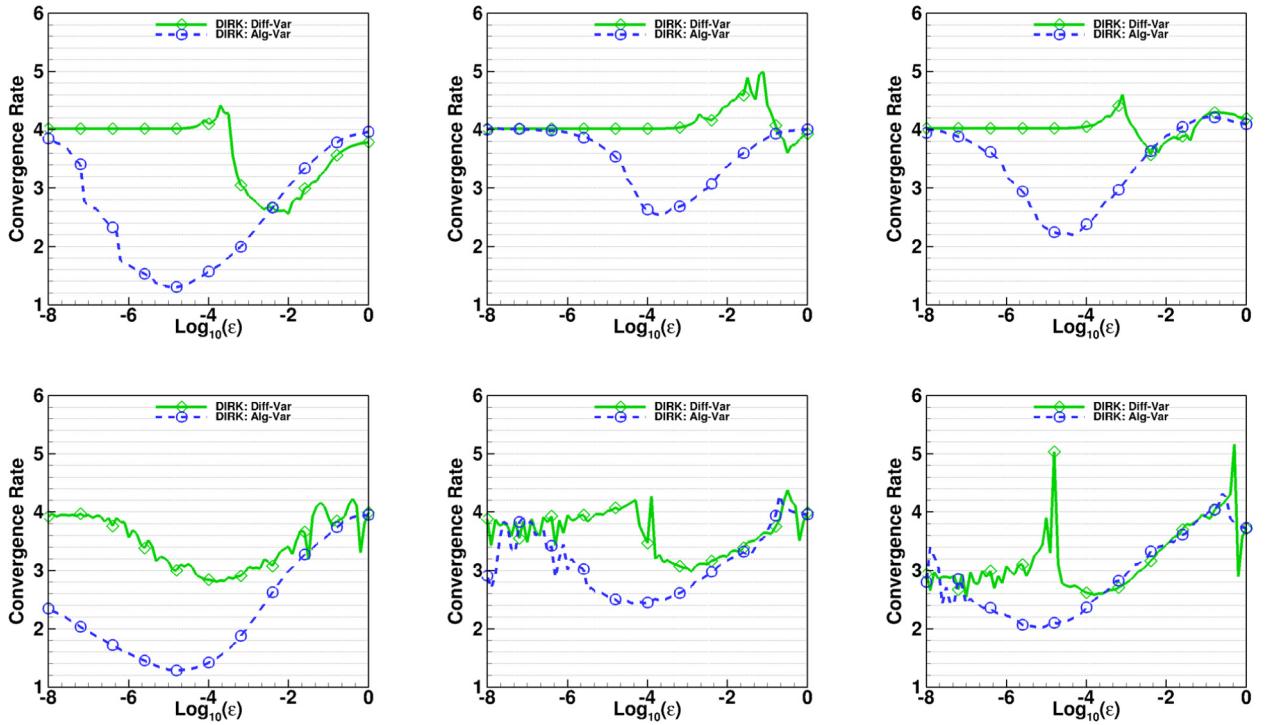


Fig. 4. A convergence rate comparison between SDIRK4 (left), ESDIRK4(3)6L[2]SA₁ (middle) and ESDIRK4(3)7L[2]SA₁ (right) on Kaps' problem (top) and van der Pol's equation (bottom).

was done by simply reducing γ from 0.250 to 0.248. This resulted in a slight improvement in most method properties but there is likely to be little discernible difference between the methods across various singular perturbation (or real-world) problems. Hence results for ESDIRK4(3)6L[2]SA₂ will not be plotted below. The three methods considered are then: SDIRK4, ESDIRK4(3)6L[2]SA₁ and ESDIRK4(3)7L[2]SA. SDIRK4 has been a very popular method for decades [9]. Its weakness is that it has a stage-order of only one. In [14], the stage-order two method, ESDIRK4(3)6L[2]SA₁, was recommended as a good general-purpose method. By adding a seventh-stage, as is done in ESDIRK4(3)7L[2]SA, the important quantity γ can be cut in half. This can be expected to have a positive effect on the iterative convergence behavior of the method. Besides improving iterative convergence, the method is substantially more accurate (formally) than ESDIRK4(3)6L[2]SA₂. Figs. 4, 5 and 6 show the convergence rates and error-versus-work performance of these three methods on the two test problems.

Fig. 4 shows differences in convergence rates between SDIRK4, ESDIRK4(3)6L[2]SA₁ and ESDIRK4(3)7L[2]SA₁. Though the plots are less than smooth, it may be clearly seen that SDIRK4 suffers from severe order reduction of the emerging algebraic variable on both test problems. The stage-order two methods exhibit much better convergence behavior on the algebraic variable. Focusing on the error-versus-work performance of the fourth-order methods, comparisons may be made at different values of the stiffness parameter with Kaps' problem. These are shown in Fig. 5 for the differential and algebraic variables. While the differences in the differential variable between methods is relatively small, SDIRK4 has much more difficulty on the algebraic variable at stiffnesses of $\epsilon = 10^{-3}$ and $\epsilon = 10^{-6}$. The six- and seven-stage methods behave similarly. Repeating the error-versus-work performance study using the more demanding van der Pol equation rather than Kaps' problem, Fig. 6 shows differences in convergence rates between the same three methods. The behavior of all three methods appears to be significantly more unsettled than with Kaps's problem. Even the accuracy of the differential variable is degraded at intermediate stiffnesses with van der Pol's equation. Once again, the performance of the algebraic variable with SDIRK4 is worse than with the other methods. What is intriguing is the analogy between methods with the minimal stages for a given order and those with an extra stage. At third-order, this was ESDIRK3(2)4L[2]SA and ESDIRK3(2)5L[2]SA. At fourth-order, it is ESDIRK4(3)6L[2]SA₁ and ESDIRK4(3)7L[2]SA. By adding an extra stage at orders three and four, γ was approximately cut in half and leading order errors were reduced substantially. However, on a problem as severe as van der Pol's problem, it appears that increasing the number of stages may degrade the performance, in spite of the seeming advantages of the increased stage. Part of this is likely the fact that performance is related to work and work is measured in the number of implicit stages of the method.

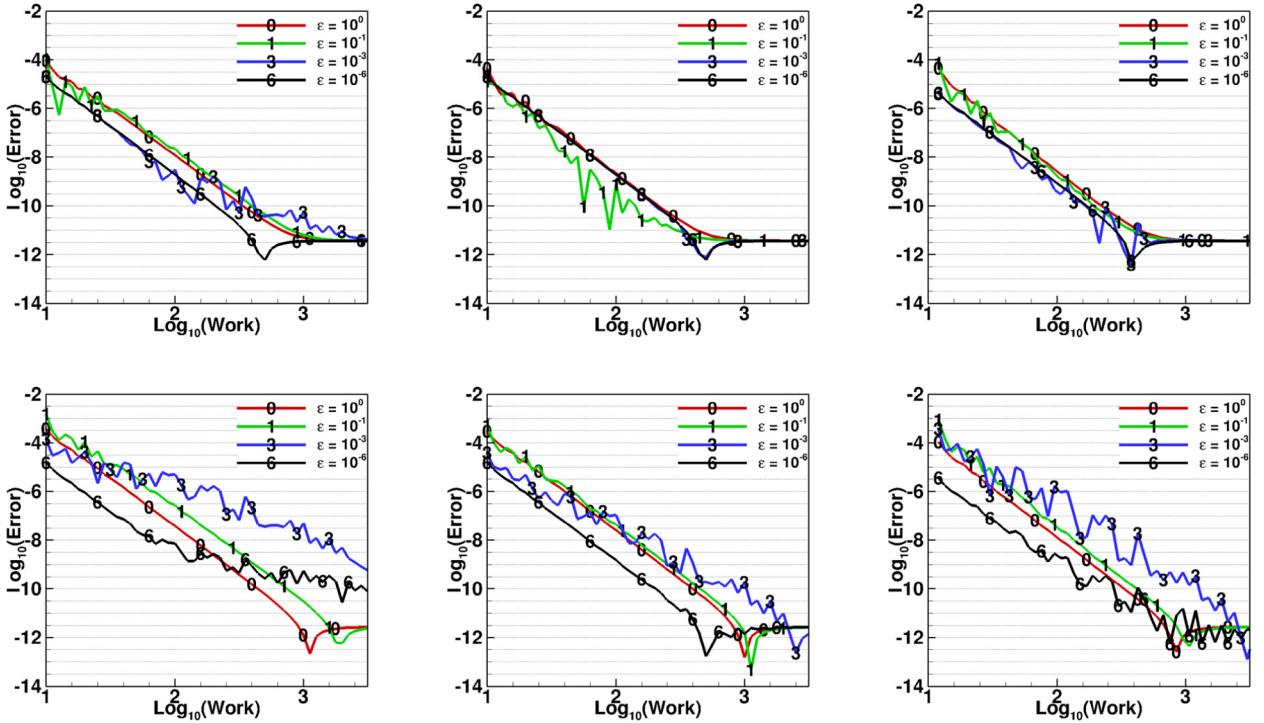


Fig. 5. An error versus work comparison between SDIRK4 (left), ESDIRK4(3)6L[2]SA₁ (middle) and ESDIRK4(3)7L[2]SA₁ (right) for the differential (top) and algebraic (bottom) variable on Kaps' problem.

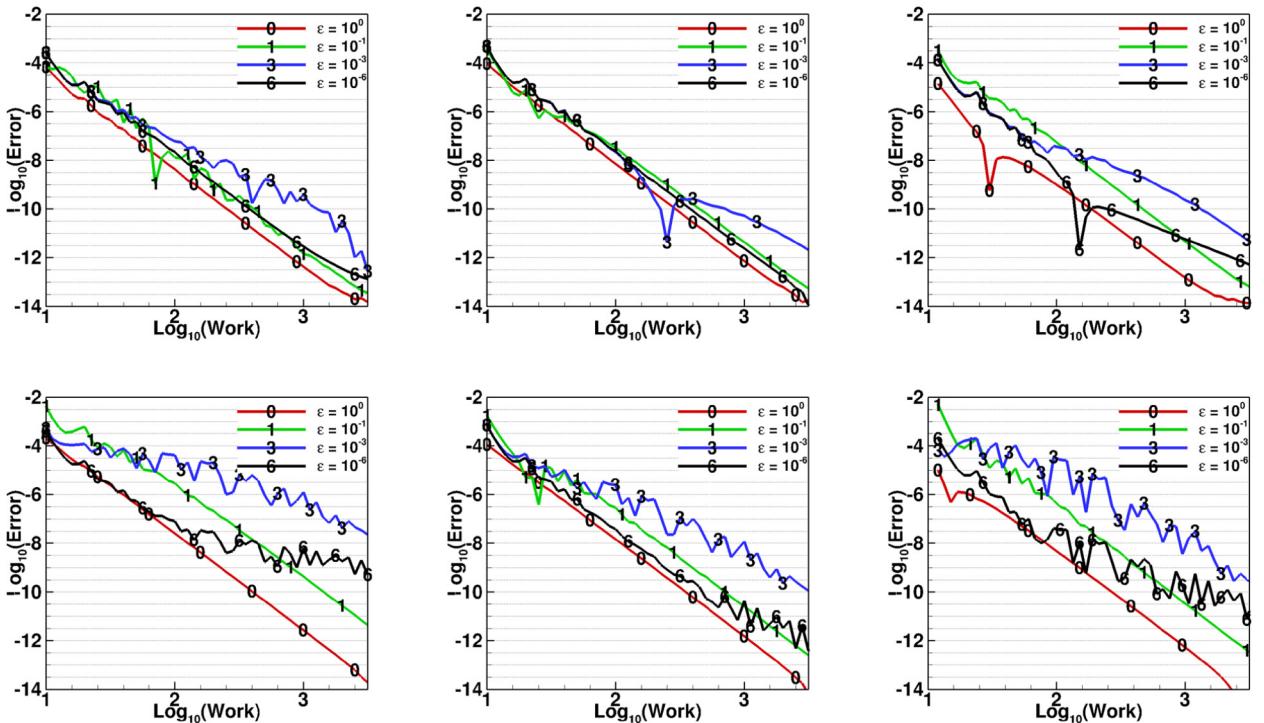


Fig. 6. An error versus work comparison between SDIRK4 (left), ESDIRK4(3)6L[2]SA₁ (middle) and ESDIRK4(3)7L[2]SA₁ (right) for the differential (top) and algebraic (bottom) variable on van der Pol's equation.

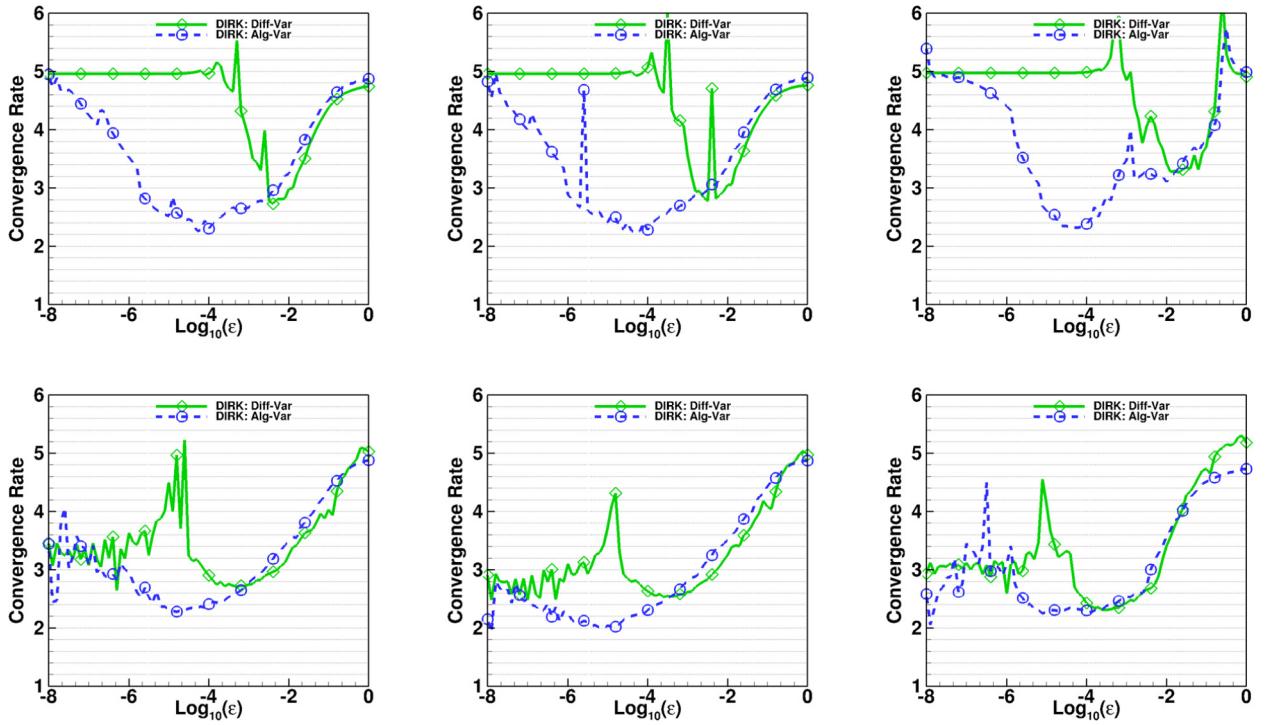


Fig. 7. A convergence rate comparison between ESDIRK5(4)7L[2]SA₁ (left), ESDIRK5(4)7L[2]SA₂ (middle) and ESDIRK5(4)8L[2]SA₁ (right) on Kaps' problem (top) and van der Pol's equation (bottom).

12.1.3. Fifth-order methods

Convergence rate testing of fifth-order ESDIRKs involved three methods: ESDIRK5(4)7L[2]SA₁, ESDIRK5(4)7L[2]SA₂ and ESDIRK5(4)8L[2]SA. Method ESDIRK5(4)7L[2]SA₁ was given in [14] but did not behave as well as had been hoped. This was likely caused by either the lack of L-stability or strong nonlinear instability on stage six. To address these issues, ESDIRK5(4)7L[2]SA₂ was created. Unlike its predecessor, it is internally L-stable on all stages and is not significantly nonlinearly unstable on any stage. Further, the leading order error was reduced by 30%. In attempt to see if adding an additional stage could improve this class of methods, ESDIRK5(4)8L[2]SA was created. It is potentially slightly more efficient than ESDIRK5(4)7L[2]SA₂ by virtue of having a reduced value of the product $\gamma(s-1)$. It also offers an L-stable embedded method and a reduced leading-order error relative to ESDIRK5(4)7L[2]SA₂. However, with strong order-reduction present, lower leading-order error is not significant.

Fig. 7 shows differences in convergence rates between ESDIRK5(4)7L[2]SA₁, ESDIRK5(4)7L[2]SA₂ and ESDIRK5(4)8L[2]SA₁. Within the irregular nature of the graphs, there is little to distinguish between the methods. Alternatively, an error-versus-work comparison may be made between these same four methods at different values of the stiffness parameter using Kaps' problem. These are shown in Fig. 8 for the differential and algebraic variables. For the differential variable, the eight-stage method appears to be slightly more efficient. It may also be more efficient for the algebraic variable.

Conducting an error-versus-work comparison with these same fifth-order methods using van der Pol's equation may also be done. Fig. 9 shows differences in convergence rates between ESDIRK5(4)7L[2]SA₁, ESDIRK5(4)7L[2]SA₂ and ESDIRK5(4)8L[2]SA₁. These are shown in for the differential and algebraic variables. All three methods behave similarly with modest order reduction, particularly with the algebraic variable. Rather than looking for differences in performance, it might be concluded that all three methods are behaving as a high-quality, fifth-order methods should when it is stiffly-accurate, is stage-order two and with most of its internal stages being L-stable. To distinguish between these methods, other criteria are needed. Both ESDIRK5(4)7L[2]SA₂ and ESDIRK5(4)8L[2]SA would seem to be the best choices at fifth-order given their improved internal stability relative to ESDIRK5(4)7L[2]SA₁.

12.1.4. Sixth-order method

The left side of Fig. 10 shows differences in convergence rates of the differential and emerging algebraic variable for ESDIRK6(5)9L[2]SA on Kaps' and van der Pol's Equation. Alternatively, an error versus work comparison may be made with this method at different values of the stiffness parameter. These are shown in Fig. 10 for the differential (middle)

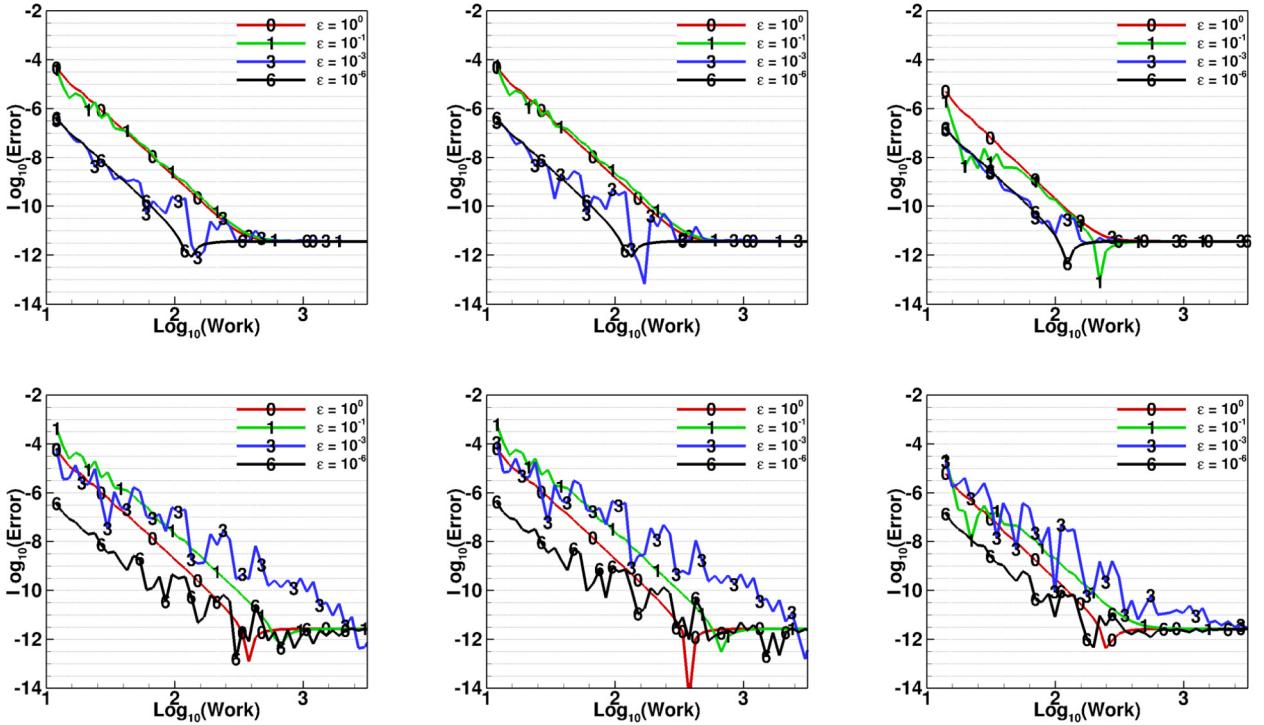


Fig. 8. An error versus work comparison between ESDIRK5(4)7L[2]SA₁ (left), ESDIRK5(4)7L[2]SA₂ (middle) and ESDIRK5(4)8L[2]SA₁ (right) for the differential (top) and algebraic (bottom) variables on Kaps' problem.

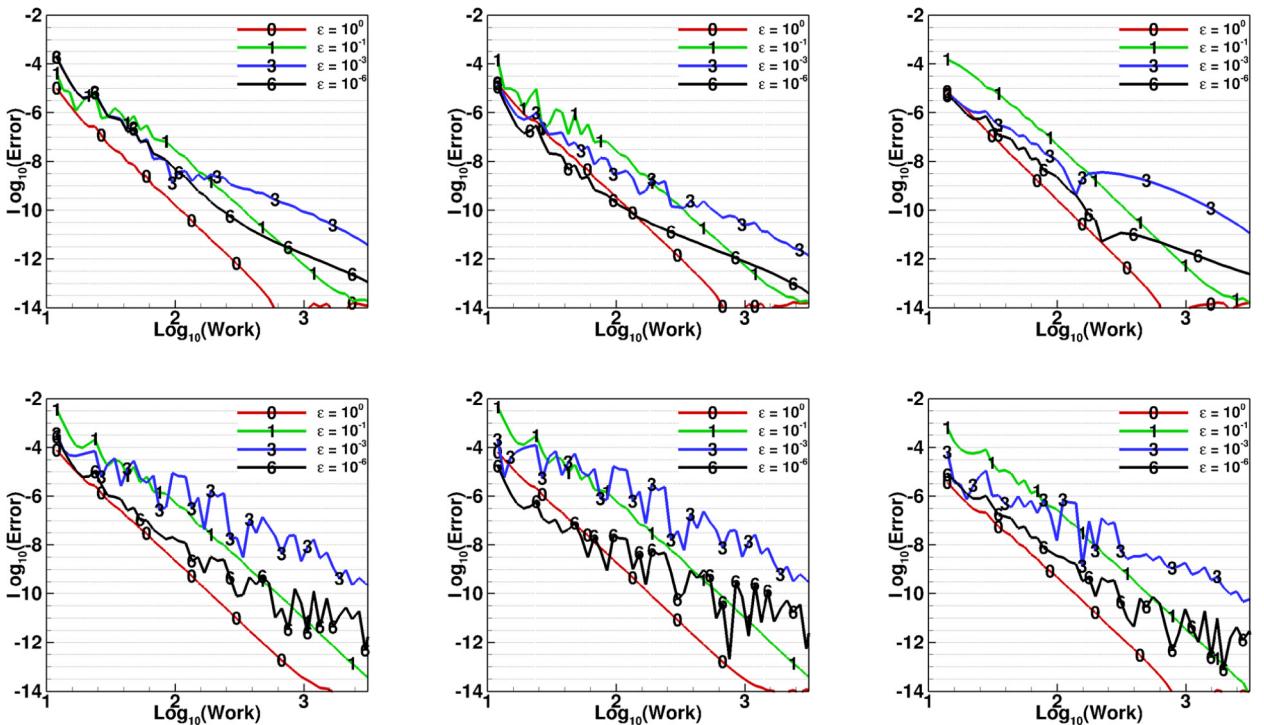


Fig. 9. An error versus work comparison between ESDIRK5(4)7L[2]SA₁ (left), ESDIRK5(4)7L[2]SA₂ (middle) and ESDIRK5(4)8L[2]SA₁ (right) for the differential (top) and algebraic (bottom) variables on van der Pol Equation.

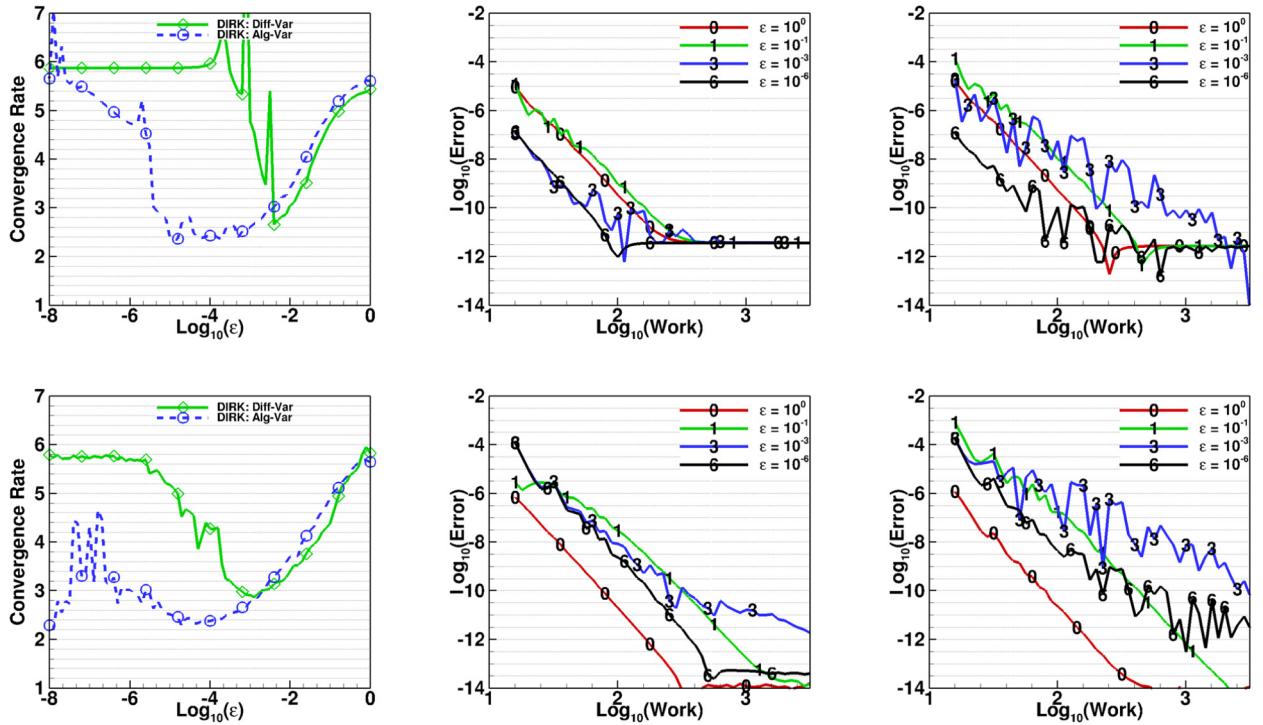


Fig. 10. A convergence rate (left) and error versus work comparison between ESDIRK6(5)9L[2]SA for the differential (middle) and algebraic (right) variables on Kaps' problem (top) and van der Pol's equation (bottom).

and algebraic (right) variables, for both test equations. Order reduction is significant for this method, as should be expected.

13. Conclusions

This paper is based on a comprehensive review of DIRK-type methods [14] applied to first-order ODEs but it extends the review with five new general-purpose ESDIRK methods based on lessons learned. Though it does not include a dense-output method nor stage-value predictors for these new methods, the methods were designed with these accoutrements in mind.

Based on the review of method characteristics, these methods focus on having a stage-order of two, stiff-accuracy, L-stability, internal L-stability, a high quality embedded method, small magnitudes of the algebraic stability matrix eigenvalues and small values of a_{ii} . These choices are also consistent with maximizing scheme efficiency. As stage-order governs the severity of order reduction, focusing on stage-order two methods facilitates accuracy.

Observed order reduction is problem dependent. Methods exhibited moderate order reduction on Kaps' problem but more severe order reduction on the vdP problem. Therefore, the choice of the optimal method is to some extent a function of the problem severity. For problems exhibiting extremely stiff behavior, the utility of the fifth- and sixth-order methods presented in this paper is likely lower than that of the lower-order methods.

At third-order, we prefer ESDIRK3(2)5L[2]SA while at fourth-order, ESDIRK4(3)6L[2]SA₁ ESDIRK4(3)6L[2]SA₂ and ESDIRK4(3)7L[2]SA are all good choices. Provided the problem at hand is not exceedingly stiff and the required error tolerance is tight, methods ESDIRK5(4)7L[2]SA₂ and ESDIRK5(4)8L[2]SA are recommended. Only in cases with quite modest stiffness does it make sense to use ESDIRK6(5)9L[2]SA. Any of these choices could be combined with the H321 or PPID error-controllers shown in Table 5.

Acknowledgements

The second author's funding was provided by Dr. Mujeeb R. Malik, program manager of NASA's Transformational Tools and Technologies (TTT) project. Both authors would like to thank the reviewer for their recommendations that have improved this paper.

- [10] E. Hairer, S.P. Nørsett, G. Wanner, Solving Ordinary Differential Equations, I: Nonstiff Problems, 2 ed., Springer-Verlag, Berlin, 1993.
- [11] I. Higueras, T. Roldán, Order barrier for low-storage DIRK methods with positive weights, *J. Sci. Comput.* 75 (1) (2018) 395–404.
- [12] G. Izzo, Z. Jackiewicz, Highly stable implicit-explicit Runge–Kutta methods, *Appl. Numer. Math.* 113 (2017) 71–92.
- [13] C.A. Kennedy, M.H. Carpenter, Additive Runge–Kutta schemes for convection-diffusion-reaction equations, *Appl. Numer. Math.* 44 (1–2) (2003) 139–181.
- [14] C.A. Kennedy, M.H. Carpenter, Diagonally Implicit Runge–Kutta Methods for Ordinary Differential Equations: A Review, NASA/TM-2016-219173, NASA Langley Research Center, 2016, 162 pp.
- [15] C.A. Kennedy, M.H. Carpenter, Higher-order additive Runge–Kutta schemes for ordinary differential equations, *Appl. Numer. Math.* 136 (2019) 183–205.
- [16] P.J. Prince, J.R. Dormand, High order embedded Runge–Kutta formulae, *J. Comput. Appl. Math.* 7 (1) (1981) 67–75.